

# On the Effect of Bias Estimation on Coverage Accuracy in Nonparametric Inference

Sebastian Calonico	Matias D. Cattaneo
Department of Economics	Department of Economics
University of Miami	Department of Statistics
Coral Gables, FL 33124	University of Michigan
	Ann Arbor, MI 48109

Max H. Farrell  
Booth School of Business  
University of Chicago  
Chicago, IL 60637

November 11, 2016

---

Sebastian Calonico is Assistant Professor of Economics, Department of Economics, University of Miami, Coral Gables, FL 33124 (email: scalonico@bus.miami.edu). Matias D. Cattaneo is Associate Professor of Economics and Statistics, Department of Economics and Department of Statistics, University of Michigan, Ann Arbor, MI 48109 (email: cattaneo@umich.edu). Max H. Farrell is Assistant Professor of Econometrics and Statistics, Booth School of Business, University of Chicago, Chicago, IL 60637 (email: max.farrell@chicagobooth.edu). The second author gratefully acknowledges financial support from the National Science Foundation (SES 1357561 and SES 1459931). We thank Ivan Canay, Xu Cheng, Joachim Freyberger, Bruce Hansen, Joel Horowitz, Michael Jansson, Francesca Molinari, Ulrich Müller, and Andres Santos for thoughtful comments and suggestions, as well as seminar participants at Cornell, Cowles Foundation, CREST Statistics, London School of Economics, Northwestern, Ohio State University, Princeton, Toulouse School of Economics, University of Bristol, and University College London. The Associate Editor and three reviewers also provided very insightful comments that improved this manuscript.

## Abstract

Nonparametric methods play a central role in modern empirical work. While they provide inference procedures that are more robust to parametric misspecification bias, they may be quite sensitive to tuning parameter choices. We study the effects of bias correction on confidence interval coverage in the context of kernel density and local polynomial regression estimation, and prove that bias correction can be preferred to undersmoothing for minimizing coverage error and increasing robustness to tuning parameter choice. This is achieved using a novel, yet simple, Studentization, which leads to a new way of constructing kernel-based bias-corrected confidence intervals. In addition, for practical cases, we derive coverage error optimal bandwidths and discuss easy-to-implement bandwidth selectors. For interior points, we show that the MSE-optimal bandwidth for the original point estimator (before bias correction) delivers the fastest coverage error decay rate after bias correction when second-order (equivalent) kernels are employed, but is otherwise suboptimal because it is too “large”. Finally, for odd-degree local polynomial regression, we show that, as with point estimation, coverage error adapts to boundary points automatically when appropriate Studentization is used; however, the MSE-optimal bandwidth for the original point estimator is suboptimal. All the results are established using valid Edgeworth expansions and illustrated with simulated data. Our findings have important consequences for empirical work as they indicate that bias-corrected confidence intervals, coupled with appropriate standard errors, have smaller coverage error and are less sensitive to tuning parameter choices in practically relevant cases where additional smoothness is available.

**Keywords:** Edgeworth expansion, coverage error, kernel methods, local polynomial regression.

# 1 Introduction

Nonparametric methods are widely employed in empirical work, as they provide point estimates and inference procedures that are more robust to parametric misspecification bias. Kernel-based methods are commonly used to estimate densities, conditional expectations, and related functions nonparametrically in a wide variety of settings. However, these methods require specifying a bandwidth and their performance in applications crucially relies on how this tuning parameter is chosen. In particular, valid inference requires the delicate balancing act of selecting a bandwidth small enough to remove smoothing bias, yet large enough to ensure adequate precision. Tipping the scale in either direction can greatly skew results. This paper studies kernel density and local polynomial regression estimation and inference based on the popular Wald-type statistics and demonstrates (via higher-order expansions) that by coupling explicit bias correction with a novel, yet simple, Studentization, inference can be made substantially more robust to bandwidth choice, greatly easing implementability.

Perhaps the most common bandwidth selection approach is to minimize the asymptotic mean-square error (MSE) of the point estimator, and then use this bandwidth choice even when the goal is inference. So difficult is bandwidth selection perceived to be, that despite the fact that the MSE-optimal bandwidth leads to *invalid* confidence intervals, even asymptotically, this method is still advocated, and is the default in most popular software. Indeed, [Hall and Kang \(2001, p. 1446\)](#) write: “there is a growing belief that the most appropriate approach to constructing confidence regions is to estimate [the density] in a way that is optimal for pointwise accuracy. . . . [I]t has been argued that such an approach has advantages of clarity, simplicity and easy interpretation.”

The underlying issue, as formalized below, is that bias must be removed for valid inference, and (in particular) the MSE-optimal bandwidth is “too large”, leaving a bias that is still first order. Two main methods have been proposed to address this: undersmoothing and explicit bias correction. We seek to compare these two, and offer concrete ways to better implement the latter. Undersmoothing amounts to choosing a bandwidth smaller than would be optimal

for point estimation, then argue that the bias is smaller than the variability of the estimator asymptotically, leading to valid distributional approximations and confidence intervals. In practice this method often involves simply shrinking the MSE-optimal bandwidth by an ad-hoc amount. The second approach is to bias correct the estimator with the explicit goal of removing the bias that caused the invalidity of the inference procedure in the first place.

It has long been believed that undersmoothing is preferable for two reasons. First, theoretical studies showed inferior asymptotic coverage properties of bias-corrected confidence intervals. The pivotal work was done by [Hall \(1992b\)](#), and has been relied upon since. Second, implementation of bias correction is perceived as more complex because a second (usually different) bandwidth is required, deterring practitioners. However, we show theoretically that bias correction is always as good as undersmoothing, and better in many practically relevant cases, if the new standard errors that we derive are used. Further, our findings have important implications for empirical work because the resulting confidence intervals are more robust to bandwidth choice, including to the bandwidth used for bias estimation. Indeed, the two bandwidths can be set equal, a simple and automatic choice that performs well in practice and is even optimal in certain objective senses.

Our proposed robust bias correction method delivers valid confidence intervals (and related inference procedures) even when using the MSE-optimal bandwidth for the original point estimator, the most popular approach in practice. Moreover, we show that at interior points, when using second-order kernels or local linear regressions, the coverage error of such intervals vanishes at the best possible rate. (Throughout, the notion of “optimal” or “best” rate is defined as the fastest achievable coverage decay for a *fixed* kernel order or polynomial degree; and is also different from optimizing point estimation.) When higher-order kernels are used, or boundary points are considered, we find that the corresponding MSE-optimal bandwidth leads to asymptotically valid intervals, but with suboptimal coverage error rates, and must be shrunk (sometimes considerably) for better inference.

Heuristically, employing the MSE-optimal bandwidth for the original point estimator, prior

to bias correction, is like undersmoothing the bias-corrected point estimator, though the latter estimator employs a possibly random,  $n$ -varying kernel, and requires a different Studentization scheme. It follows that the conventional MSE-optimal bandwidth commonly used in practice need not be optimal, even after robust bias correction, when the goal is inference. Thus, we present new coverage error optimal bandwidths and a fully data-driven direct plug-in implementation thereof, for use in applications. In addition, we study the important related issue of asymptotic length of the new confidence intervals.

Our comparisons of undersmoothing and bias correction are based on Edgeworth expansions for density estimation and local polynomial regression, allowing for different levels of smoothness of the unknown functions. We prove that explicit bias correction, coupled with our proposed standard errors, yields confidence intervals with coverage that is as accurate, or better, than undersmoothing (or, equivalently, yields dual hypothesis tests with lower error in rejection probability). Loosely speaking, this improvement is possible because explicit bias correction can remove more bias than undersmoothing, while our proposed standard errors capture not only the variability of the original estimator but also the additional variability from bias correction. To be more specific, our robust bias correction approach yields higher-order refinements whenever additional smoothness is available, and is asymptotically equivalent to the best undersmoothing procedure when no additional smoothness is available.

Our findings contrast with well established recommendations: [Hall \(1992b\)](#) used Edgeworth expansions to show that undersmoothing produces more accurate intervals than explicit bias correction in the density case and [Neumann \(1997\)](#) repeated this finding for kernel regression. The key distinction is that their expansions, while imposing the same levels of smoothness as we do, crucially relied on the assumption that the bias correction was first-order negligible, essentially forcing bias correction to remove less bias than undersmoothing. In contrast, we allow the bias estimator to potentially have a first order impact, an alternative asymptotic experiment designed to more closely mimic the finite-sample behavior of bias correction. Therefore, our results formally show that whenever additional smoothness is

available to characterize leading bias terms, as it is usually the case in practice where MSE-optimal bandwidth are employed, our robust bias correction approach yields higher-order improvements relative to standard undersmoothing.

Our standard error formulas are based on fixed- $n$  calculations, as opposed to asymptotics, which also turns out to be important. We show that using asymptotic variance formulas can introduce further errors in coverage probability, with particularly negative consequences at boundary points. This turns out to be at the heart of the “quite unexpected” conclusion found by [Chen and Qin \(2002, Abstract\)](#) that local polynomial based confidence intervals are not boundary-adaptive in coverage error: we prove that this is not the case with proper Studentization. Thus, as a by-product of our main theoretical work, we establish higher-order boundary carpentry of local polynomial based confidence intervals that use a fixed- $n$  standard error formula, a result that is of independent (but related) interest.

This paper is connected to the well-established literature on nonparametric smoothing, see [Wand and Jones \(1995\)](#), [Fan and Gijbels \(1996\)](#), [Horowitz \(2009\)](#), and [Ruppert et al. \(2009\)](#) for reviews. For more recent work on bias and related issues in nonparametric inference, see [Hall and Horowitz \(2013\)](#), [Calonico et al. \(2014\)](#), [Hansen \(2015\)](#), [Armstrong and Kolesár \(2015\)](#), [Schennach \(2015\)](#), and references therein. We also contribute to the literature on Edgeworth expansions, which have been used both in parametric and, less frequently, nonparametric contexts; see, e.g., [Bhattacharya and Rao \(1976\)](#) and [Hall \(1992a\)](#). Fixed- $n$  versus asymptotic-based Studentization has also captured some recent interest in other contexts, e.g., [Mykland and Zhang \(2015\)](#). Finally, see [Calonico et al. \(2016\)](#) for uniformly valid Edgeworth expansions and optimal inference in the context of regression discontinuity designs.

The paper proceeds as follows. Section 2 studies density estimation at interior points and states the main results on error in coverage probability and its relationship to bias reduction and underlying smoothness, as well as discussing bandwidth choice and interval length. Section 3 then studies local polynomial estimation, at interior and boundary points. Practical guidance is explicitly discussed in Sections 2.4 and 3.3, respectively; all methods are available in the R

package `nprobust` on CRAN. Section 4 summarizes the results of a Monte Carlo study, and Section 5 concludes. Some technical details, all proofs, and additional simulation evidence are collected in a lengthy online supplement.

## 2 Density Estimation and Inference

We first present our main ideas and conclusions for inference on the density at an interior point, as this requires relatively little notation. The data are assumed to obey the following.

**Assumption 2.1** (Data-generating process).  *$\{X_1, \dots, X_n\}$  is a random sample with an absolutely continuous distribution with Lebesgue density  $f$ . In a neighborhood of  $x$ ,  $f > 0$ ,  $f$  is  $S$ -times continuously differentiable with bounded derivatives  $f^{(s)}$ ,  $s = 1, 2, \dots, S$ , and  $f^{(S)}$  is Hölder continuous with exponent  $\varsigma$ .*

The parameter of interest is  $f(x)$  for a fixed scalar point  $x$  in the interior of the support. [In the supplemental appendix we discuss how our results extend naturally to multivariate  $X_i$  and derivative estimation.] The classical kernel-based estimator of  $f(x)$  is

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right), \quad (1)$$

for a kernel function  $K$  that integrates to 1 and positive bandwidth  $h \rightarrow 0$  as  $n \rightarrow \infty$ . The choice of  $h$  can be delicate, and our work is motivated in part by the standard empirical practice of employing the MSE-optimal bandwidth choice for  $\hat{f}(x)$  when conducting inference.

In this vein, let us suppose for the moment that  $K$  is a kernel of order  $k$ , where  $k \leq S$  so that the MSE-optimal bandwidth can be characterized. The bias is then given by

$$\mathbb{E}[\hat{f}(x)] - f(x) = h^k f^{(k)}(x) \mu_{K,k} + o(h^k), \quad (2)$$

where  $f^{(\hat{k})}(x) := \partial^{\hat{k}} f(x) / \partial x^{\hat{k}}$  and  $\mu_{K,\hat{k}} = \int u^{\hat{k}} K(u) du / \hat{k}!$ . Computing the variance gives

$$(nh)\mathbb{V}[\hat{f}(x)] = \frac{1}{h} \left\{ \mathbb{E} \left[ K \left( \frac{x - X_i}{h} \right)^2 \right] - \mathbb{E} \left[ K \left( \frac{x - X_i}{h} \right) \right]^2 \right\}, \quad (3)$$

which is *non-asymptotic*:  $n$  and  $h$  are fixed in this calculation. Using other, first-order valid approximations, e.g.  $(nh)\mathbb{V}[\hat{f}(x)] \approx f(x) \int K(u)^2 du$ , will have finite sample consequences that manifest as additional terms in the Edgeworth expansions. In fact, Section 3 shows that using an asymptotic variance for local polynomial regression removes automatic coverage-error boundary adaptivity.

Together, the prior two displays are used to characterize the MSE-optimal bandwidth,  $h_{\text{mse}}^* \propto n^{-1/(1+2\hat{k})}$ , however, using this bandwidth leaves a bias that is too large, relative to the variance, to conduct valid inference for  $f(x)$ . To address this important practical problem, researchers must either undersmooth the point estimator (i.e., construct  $\hat{f}(x)$  with a bandwidth smaller than  $h_{\text{mse}}^*$ ) or bias-correct the point estimator (i.e., estimate and subtract the leading bias after using  $h_{\text{mse}}^*$ , or a “larger” bandwidth, to construct  $\hat{f}(x)$ ). Thus, the question we seek to answer is this: if the bias is given by (2), is one better off estimating the leading bias (explicit bias correction) or choosing  $h$  small enough to render the bias negligible (undersmoothing) when forming nonparametric confidence intervals?

To answer this question, and to motivate our new robust approach, we first detail the bias correction and variance estimators. Explicit bias correction estimates the leading term of Eqn. (2), denoted by  $B_f$ , using a kernel estimator of  $f^{(\hat{k})}(x)$ , defined as:

$$\hat{B}_f = h^{\hat{k}} \hat{f}^{(\hat{k})}(x) \mu_{K,\hat{k}}, \quad \text{where} \quad \hat{f}^{(\hat{k})}(x) = \frac{1}{nb^{1+\hat{k}}} \sum_{i=1}^n L^{(\hat{k})} \left( \frac{x - X_i}{b} \right),$$

for a kernel  $L(\cdot)$  of order  $\ell$  and a bandwidth  $b \rightarrow 0$  as  $n \rightarrow \infty$ . Importantly,  $\hat{B}_f$  takes this form for any  $\hat{k}$  and  $S$ , even if (2) fails; see Sections 2.2 and 2.3 for discussion. Conventional



Studentized statistics based on undersmoothing and explicit bias correction are, respectively,

$$T_{\text{us}}(x) = \frac{\sqrt{nh}(\hat{f}(x) - f(x))}{\hat{\sigma}_{\text{us}}} \quad \text{and} \quad T_{\text{bc}}(x) = \frac{\sqrt{nh}(\hat{f}(x) - \hat{B}_f - f(x))}{\hat{\sigma}_{\text{us}}},$$

where  $\hat{\sigma}_{\text{us}}^2 := \hat{\mathbb{V}}[\hat{f}(x)]$  is the natural estimator of the variance of  $\hat{f}(x)$  which only replaces the two expectations in (3) with sample averages, thus maintaining the nonasymptotic spirit. These are the two statistics compared in the influential paper of Hall (1992b), under the same assumption imposed herein.

From the form of these statistics, two points are already clear. First, the numerator of  $T_{\text{us}}$  relies on choosing  $h$  vanishing fast enough so that the bias is asymptotically negligible after scaling, whereas  $T_{\text{bc}}$  allows for slower decay by virtue of the manual estimation of the leading bias. Second,  $T_{\text{bc}}$  requires that the variance of  $h^{\hat{k}} \hat{f}^{(\hat{k})}(x) \mu_{K, \hat{k}}$  be first-order asymptotically negligible:  $\hat{\sigma}_{\text{us}}$  in the denominator only accounts for the variance of the main estimate, but  $\hat{f}^{(\hat{k})}(x)$ , being a kernel-based estimator, naturally has a variance controlled by its bandwidth. That is, even though  $\hat{\sigma}_{\text{us}}^2$  is based on a fixed- $n$  calculation, the variance of the numerator of  $T_{\text{bc}}$  only coincides with the denominator asymptotically. Under this regime, Hall (1992b) showed that the bias reduction achieved in  $T_{\text{bc}}$  is too expensive in terms of noise and that undersmoothing dominates explicit bias correction for coverage error.

We argue that there need not be such a “mismatch” between the numerator of the bias-corrected statistic and the Studentization, and thus consider a third option corresponding to the idea of capturing the finite sample variability of  $\hat{f}^{(\hat{k})}(x)$  directly. To do so, note that we may write, after setting  $\rho = h/b$ ,

$$\hat{f}(x) - h^{\hat{k}} \hat{f}^{(\hat{k})}(x) \mu_{K, \hat{k}} = \frac{1}{nh} \sum_{i=1}^n M\left(\frac{x - X_i}{h}\right), \quad M(u) = K(u) - \rho^{1+\hat{k}} L^{(\hat{k})}(\rho u) \mu_{K, \hat{k}}. \quad (4)$$

We then define the collective variance of the density estimate and the bias correction as  $\sigma_{\text{rbc}}^2 = (nh) \mathbb{V}[\hat{f}(x) - \hat{B}_f]$ , exactly as in Eqn. (3), but with  $M(\cdot)$  in place of  $K(\cdot)$ , and its estimator  $\hat{\sigma}_{\text{rbc}}^2$  exactly as  $\hat{\sigma}_{\text{us}}^2$ . Therefore, our proposed robust bias corrected inference approach

is based on

$$T_{\text{rbc}} = \frac{\sqrt{nh}(\hat{f}(x) - h^{\hat{k}} \hat{f}^{(\hat{k})}(x) \mu_{K,\hat{k}} - f(x))}{\hat{\sigma}_{\text{rbc}}}.$$

That is, our proposed standard errors are based on a fixed- $n$  calculation that captures the variability of both  $\hat{f}(x)$  and  $\hat{f}^{(\hat{k})}(x)$ , and their covariance. As shown in Section 3, the case of local polynomial regression is qualitatively analogous, but notationally more complicated.

The quantity  $\rho = h/b$  is key. If  $\rho \rightarrow 0$ , then the second term of  $M$  is dominated by the first, i.e. the bias correction is first-order negligible. In this case,  $\sigma_{\text{us}}^2$  and  $\sigma_{\text{rbc}}^2$  (and their estimators) will be first-order, but not higher-order, equivalent. This is exactly the sense in which traditional bias correction relies on an asymptotic variance, instead of a fixed- $n$  one, and pays the price in coverage error: for any finite sample, standard bias corrected inference can substantially impact the results. To more accurately capture finite sample behavior of bias correction we allow  $\rho$  to converge to any (nonnegative) finite limit, allowing (but not requiring) the bias correction to be first-order important, unlike prior work. We show that doing so yields more accurate confidence intervals (i.e., higher-order corrections).

## 2.1 Generic Higher Order Expansions of Coverage Error

We first present generic Edgeworth expansions for all three procedures (undersmoothing, traditional bias correction and robust bias correction), which are agnostic regarding the level of available smoothness (controlled by  $S$  in Assumption 2.1). To be specific, we give higher-order expansions of the error in coverage probability of the following  $(1 - \alpha)\%$  confidence intervals based on Normal approximations for the statistics  $T_{\text{us}}$ ,  $T_{\text{bc}}$ , and  $T_{\text{rbc}}$ :

$$\begin{aligned} I_{\text{us}} &= \left[ \hat{f} - z_{1-\frac{\alpha}{2}} \frac{\hat{\sigma}_{\text{us}}}{\sqrt{nh}}, \hat{f} - z_{\frac{\alpha}{2}} \frac{\hat{\sigma}_{\text{us}}}{\sqrt{nh}} \right], \\ I_{\text{bc}} &= \left[ \hat{f} - \hat{B}_f - z_{1-\frac{\alpha}{2}} \frac{\hat{\sigma}_{\text{us}}}{\sqrt{nh}}, \hat{f} - \hat{B}_f - z_{\frac{\alpha}{2}} \frac{\hat{\sigma}_{\text{us}}}{\sqrt{nh}} \right], \quad \text{and} \\ I_{\text{rbc}} &= \left[ \hat{f} - \hat{B}_f - z_{1-\frac{\alpha}{2}} \frac{\hat{\sigma}_{\text{rbc}}}{\sqrt{nh}}, \hat{f} - \hat{B}_f - z_{\frac{\alpha}{2}} \frac{\hat{\sigma}_{\text{rbc}}}{\sqrt{nh}} \right], \end{aligned} \tag{5}$$

where  $z_\alpha$  is the upper  $\alpha$ -percentile of the Gaussian distribution. Here and in the sequel we omit the point of evaluation  $x$  for simplicity. Equivalently, our results can characterize the error in rejection probability of the corresponding hypothesis tests. In the following subsections, we give specific results under different smoothness assumptions and make direct comparisons of the methods.

We require the following standard conditions on the kernels  $K$  and  $L$ .

**Assumption 2.2** (Kernels). *The kernels  $K$  and  $L$  are bounded, even functions with support  $[-1, 1]$ , and are of order  $k \geq 2$  and  $\ell \geq 2$ , respectively, where  $k$  and  $\ell$  are even integers. That is,  $\mu_{K,0} = 1$ ,  $\mu_{K,k} = 0$  for  $1 \leq k < k$ , and  $\mu_{K,k} \neq 0$  and bounded, and similarly for  $\mu_{L,\ell}$  with  $\ell$  in place of  $k$ . Further,  $L$  is  $k$ -times continuously differentiable. For all integers  $k$  and  $l$  such that  $k + l = k - 1$ ,  $f^{(k)}(x_0)L^{(l)}((x_0 - x)/b) = 0$  for  $x_0$  in the boundary of the support.*

The boundary conditions are needed for the derivative estimation inherent in bias correction, even if  $x$  is an interior point, and are satisfied if the support of  $f$  is the whole real line. Higher order results also require a standard  $n$ -varying Cramér's condition, given in the supplement to conserve space (see Section S.I.3). Altogether, our assumptions are identical to those of [Hall \(1991, 1992b\)](#).

To state the results some notation is required. First, let the (scaled) biases of the density estimator and the bias-corrected estimator be  $\eta_{\text{us}} = \sqrt{nh}(\mathbb{E}[\hat{f}] - f)$  and  $\eta_{\text{bc}} = \sqrt{nh}(\mathbb{E}[\hat{f} - \hat{B}_f] - f)$ . Next, let  $\phi(z)$  be the standard Normal density, and for any kernel  $K$  define

$$\begin{aligned} q_1(K) &= \vartheta_{K,2}^{-2} \vartheta_{K,4} (z_{\frac{\alpha}{2}}^3 - 3z_{\frac{\alpha}{2}}) / 6 - \vartheta_{K,2}^{-3} \vartheta_{K,3}^2 [2z^3/3 + (z_{\frac{\alpha}{2}}^5 - 10z_{\frac{\alpha}{2}}^3 + 15z_{\frac{\alpha}{2}}) / 9], \\ q_2(K) &= -\vartheta_{K,2}^{-1} z_{\frac{\alpha}{2}}, \quad \text{and} \quad q_3(K) = \vartheta_{K,2}^{-2} \vartheta_{K,3} (2z_{\frac{\alpha}{2}}^3 / 3), \end{aligned}$$

where  $\vartheta_{K,k} = \int K(u)^k du$ . All that is conceptually important is that these functions are known, odd polynomials in  $z$  with coefficients that depend only on the kernel, and not on the sample or data generating process. Our main theoretical result for density estimation is the following.

**Theorem 1.** *Let Assumptions 2.1, 2.2, and Cramér's condition hold and  $nh/\log(nh) \rightarrow \infty$ .*

(a) If  $\eta_{\text{us}} \rightarrow 0$ , then

$$\mathbb{P}[f \in I_{\text{us}}] = 1 - \alpha + \left\{ \frac{1}{nh} q_1(K) + \eta_{\text{us}}^2 q_2(K) + \frac{\eta_{\text{us}}}{\sqrt{nh}} q_3(K) \right\} \frac{\phi(z_{\frac{\alpha}{2}})}{f} \{1 + o(1)\}.$$

(b) If  $\eta_{\text{bc}} \rightarrow 0$  and  $\rho \rightarrow 0$ , then

$$\begin{aligned} \mathbb{P}[f \in I_{\text{bc}}] = 1 - \alpha + \left\{ \frac{1}{nh} q_1(K) + \eta_{\text{bc}}^2 q_2(K) + \frac{\eta_{\text{bc}}}{\sqrt{nh}} q_3(K) \right\} \frac{\phi(z_{\frac{\alpha}{2}})}{f} \{1 + o(1)\} \\ + \rho^{1+\kappa} (\Omega_1 + \rho^\kappa \Omega_2) \phi(z_{\frac{\alpha}{2}}) z_{\frac{\alpha}{2}} \{1 + o(1)\}, \end{aligned}$$

for constants  $\Omega_1$  and  $\Omega_2$  given precisely in the supplement.

(c) If  $\eta_{\text{bc}} \rightarrow 0$  and  $\rho \rightarrow \bar{\rho} < \infty$ , then

$$\mathbb{P}[f \in I_{\text{rbc}}] = 1 - \alpha + \left\{ \frac{1}{nh} q_1(M) + \eta_{\text{bc}}^2 q_2(M) + \frac{\eta_{\text{bc}}}{\sqrt{nh}} q_3(M) \right\} \frac{\phi(z_{\frac{\alpha}{2}})}{f} \{1 + o(1)\}.$$

This result leaves the scaled biases  $\eta_{\text{us}}$  and  $\eta_{\text{bc}}$  generic, which is useful when considering different levels of smoothness  $S$ , the choices of  $\kappa$  and  $\ell$ , and in comparing to local polynomial results. In the next subsection, we make these quantities more precise and compare them, paying particular attention to the role of the underlying smoothness assumed.

At present, the most visually obvious feature of this result is that all the error terms are of the same form, except for the notable presence of  $\rho^{1+\kappa}(\Omega_1 + \rho^\kappa \Omega_2)$  in part (b). These are the leading terms of  $\sigma_{\text{rbc}}^2 / \sigma_{\text{us}}^2 - 1$ , consisting of the covariance of  $\hat{f}$  and  $\hat{B}_f$  (denoted by  $\Omega_1$ ) and the variance of  $\hat{B}_f$  (denoted by  $\Omega_2$ ), and are entirely due to the “mismatch” in Studentization scheme underlying  $T_{\text{bc}}$ . [Hall \(1992b\)](#) showed how these terms prevent bias correction from performing as well as undersmoothing in terms of coverage. In essence, the potential for improved bias properties do not translate into improved inference because the variance is not well-controlled: in any finite sample,  $\hat{B}_f$  would inject variability (i.e.,  $\rho = h/b > 0$  for each  $n$ ) and thus  $\rho \rightarrow 0$  may not be a good approximation. Our new Studentization does not simply remove the leading  $\rho$  terms; the entire sequence is absent. As explained below, allowing for  $\bar{\rho} = \infty$  can not reduce bias, but will inflate variance; hence restricting to  $\bar{\rho} < \infty$  capitalizes

fully on the improvements from bias correction.

## 2.2 Coverage Error and the Role of Smoothness

Theorem 1 makes no explicit assumption about smoothness beyond the requirement that the scaled biases vanish asymptotically. The fact that the error terms in parts (a) and (c) of Theorem 1 take the same form implies that comparing coverage error amounts to comparing bias, for which the smoothness  $S$  and the kernel orders  $k$  and  $\ell$  are crucial. We now make the biases  $\eta_{\text{us}}$  and  $\eta_{\text{bc}}$  concrete and show how coverage is affected.

For  $I_{\text{us}}$ , two cases emerge: (a) enough derivatives exist to allow characterization of the MSE-optimal bandwidth ( $k \leq S$ ); and (b) no such smoothness is available ( $k > S$ ), in which case the leading term of Eqn. (2) is exactly zero and the bias depends on the unknown Hölder constant. These two cases lead to the following results.

**Corollary 1.** *Let Assumptions 2.1, 2.2, and Cramér’s condition hold and  $nh/\log(nh) \rightarrow \infty$ .*

(a) *If  $k \leq S$  and  $\sqrt{nh}h^k \rightarrow 0$ ,*

$$\mathbb{P}[f \in I_{\text{us}}] = 1 - \alpha + \left\{ \frac{1}{nh} q_1(K) + nh^{1+2k} (f^{(k)})^2 \mu_{K,k}^2 q_2(K) + h^k f^{(k)} \mu_{K,k} q_3(K) \right\} \frac{\phi(z_{\frac{\alpha}{2}})}{f} \{1 + o(1)\}.$$

(b) *If  $k > S$  and  $\sqrt{nh}h^{S+\varsigma} \rightarrow 0$ ,*

$$\mathbb{P}[f \in I_{\text{us}}] = 1 - \alpha + \frac{1}{nh} \frac{\phi(z_{\frac{\alpha}{2}})}{f} q_1(K) \{1 + o(1)\} + O(nh^{1+2(S+\varsigma)} + h^{S+\varsigma}).$$

The first result is most directly comparable to Hall (1992b, §3.4), and many other past papers, which typically take as a starting point that the MSE-optimal bandwidth can be characterized. This shows that  $T_{\text{us}}$  must be undersmoothed, in the sense the MSE-optimal bandwidth is “too large” to be valid:  $h_{\text{mse}}^* \propto n^{-1/(1+2k)}$  is not allowed. In fact,  $I_{\text{us}}(h_{\text{mse}}^*)$  asymptotically undercovers because  $T_{\text{us}}(h_{\text{mse}}^*) \rightarrow_d \mathcal{N}(1, 1)$ ; for example if  $\alpha = 0.05$ ,  $\mathbb{P}[f \in I_{\text{us}}(h_{\text{mse}}^*)] \approx 0.83$ . Instead, the optimal  $h$  for coverage error, which can be characterized and estimated, is equivalent in rates to balancing variance against bias, not squared bias as in

MSE. Part (b) shows that a faster rate of coverage error decay can be obtained by taking a sufficiently high order kernel, relative to the level of smoothness  $S$ , at the expense of feasible bandwidth selection.

Turning to robust bias correction, characterization of  $\eta_{\text{bc}}$  is more complex as it has two pieces: the second-order bias of the original point estimator, and the bias of the bias estimator itself. The former is the  $o(h^{\hat{k}})$  term of Eqn. (2) and is not the target of explicit bias correction; it depends either on higher derivatives, if they are available, or on the Hölder condition otherwise. To be precise, if  $\hat{k} \leq S - 2$ , this term is  $[h^{\hat{k}+2} + o(1)]f^{(\hat{k}+2)}\mu_{K_{\text{bc}},\hat{k}+2}$ , while otherwise is known only to be  $O(h^{S+\varsigma})$ . Importantly, the bandwidth  $b$  and order  $\ell$  do not matter here, and bias reduction beyond  $O(\min\{h^{\hat{k}+2}, h^{S+\varsigma}\})$  is not possible; there is thus little or no loss in fixing  $\ell = 2$ , which we assume from now on to simplify notation.

The bias of the bias estimator also depends on the smoothness available: if enough smoothness is available the corresponding bias term can be characterized, otherwise only its order will be known. To be specific, when smoothness is not binding ( $\hat{k} \leq S - 2$ ), arguably the most practically-relevant case, the leading term of  $\mathbb{E}[\hat{B}_f] - B_f$  will be  $h^{\hat{k}}b^2f^{(\hat{k}+2)}\mu_{K,\hat{k}}\mu_{L,2}$ . Smoothness can be exhausted in two ways, either by the point estimate itself ( $\hat{k} > S$ ) or by the bias estimation ( $S - 1 \leq \hat{k} \leq S$ ), and these two cases yield  $O(h^{\hat{k}}b^{S-\hat{k}})$  and  $O(h^{\hat{k}}b^{S+\varsigma-\hat{k}})$ , respectively, which are slightly different in how they depend on the total Hölder smoothness assumed. (Complete details are in the supplement.) Note that regardless of the value of  $\hat{k}$ , we set  $\hat{B}_f = h^{\hat{k}}\hat{f}^{(\hat{k})}\mu_{K,\hat{k}}$ , even if  $\hat{k} > S$  and  $B_f \equiv 0$ .

With these calculations for  $\eta_{\text{bc}}$ , we have the following result.

**Corollary 2.** *Let Assumptions 2.1, 2.2, and Cramér's condition hold,  $nh/\log(nh) \rightarrow \infty$ ,  $\rho \rightarrow \bar{\rho} < \infty$ , and  $\ell = 2$ .*

(a) *If  $\hat{k} \leq S - 2$  and  $\sqrt{nh}h^{\hat{k}}b^2 \rightarrow 0$ ,*

$$\begin{aligned} \mathbb{P}[f \in I_{\text{rbc}}] = 1 - \alpha + & \left\{ \frac{1}{nh}q_1(M_{\bar{\rho}}) + nh^{1+2(\hat{k}+2)}(f^{(\hat{k}+2)})^2(\mu_{K,\hat{k}+2} + \bar{\rho}^{-2}\mu_{K,\hat{k}}\mu_{L,2})^2q_2(M_{\bar{\rho}}) \right. \\ & \left. + h^{\hat{k}+2}f^{(\hat{k}+2)}(\mu_{K,\hat{k}+2} + \bar{\rho}^{-2}\mu_{K,\hat{k}}\mu_{L,2})q_3(M_{\bar{\rho}}) \right\} \frac{\phi(z_{\frac{\alpha}{2}})}{f} \{1 + o(1)\}. \end{aligned}$$

(b) If  $S - 1 \leq k \leq S$  and  $\sqrt{nh}\rho^k b^{S+\varsigma} \rightarrow 0$ ,

$$\mathbb{P}[f \in I_{\text{rbc}}] = 1 - \alpha + \frac{1}{nh} \frac{\phi(z_{\frac{\alpha}{2}})}{f} q_1(M_{\bar{\rho}}) \{1 + o(1)\} + O(nh\rho^{2k}b^{2(S+\varsigma)} + \rho^k b^{S+\varsigma}).$$

(c) If  $k > S$  and  $\sqrt{nh}(h^{S+\varsigma} \vee \rho^k b^S) \rightarrow 0$ ,

$$\mathbb{P}[f \in I_{\text{rbc}}] = 1 - \alpha + \frac{1}{nh} \frac{\phi(z_{\frac{\alpha}{2}})}{f} q_1(M_{\bar{\rho}}) \{1 + o(1)\} + O(nh(h^{S+\varsigma} \vee \rho^k b^S)^2 + (h^{S+\varsigma} \vee \rho^k b^S)).$$

Part (a) is the most empirically-relevant setting, which reflects the idea that researchers first select a kernel order, then conduct inference based on that choice, taking the unknown smoothness to be nonbinding. The most notable feature of this result, beyond the formalization of the coverage improvement, is that the coverage error terms share the same structure as those of Corollary 1, with  $k$  replaced by  $k + 2$ , and represent the same conceptual objects. By virtue of our new Studentization, the leading variance remains order  $(nh)^{-1}$  and the problematic correlation terms are absent. We explicitly discuss the advantages of robust bias correction relative to undersmoothing in the following section.

Part (a) also argues for a bounded, positive  $\rho$ . First, because bias reduction beyond  $O(h^{k+2})$  is not possible,  $\rho \rightarrow \infty$  will only inflate the variance. On the other hand,  $\bar{\rho} = 0$  requires a delicate choice of  $b$  and  $\ell > 2$ , else the second bias term dominates  $\eta_{\text{bc}}$ , and the full power of the variance correction is not exploited; that is, more bias may be removed without inflating the variance rate. Hall (1992b, p. 682) remarked that if  $\mathbb{E}[\hat{f}] - f - B_f$  is (part of) the leading bias term, then “explicit bias correction [...] is even less attractive relative to undersmoothing.” We show, on the contrary, that with our proposed Studentization, it is optimal that  $\mathbb{E}[\hat{f}] - f - B_f$  is part of the dominant bias term.

Finally, in both Corollaries above the best possible coverage error decay rate (for a given  $S$  is attained by exhausting all available smoothness. This would also yield point estimators attaining the bound of Stone (1982); robust bias correction can not evade such bounds, of course. In both Corollaries, coverage is improved relative to part (a), but the constants and

optimal bandwidths can not be quantified. For robust bias correction, Corollary 2 shows that to obtain the best rate in part (b) the unknown  $f^{(\kappa)}$  must be consistently estimated and  $\rho$  must be bounded and positive, while in part (c), bias estimation merely adds noise, but this noise is fully accounted for by our new Studentization, as long as  $\rho \rightarrow 0$  ( $b \not\rightarrow 0$  is allowed).

## 2.3 Comparing Undersmoothing and Robust Bias Correction

We now employ Corollaries 1 and 2 to directly compare nonparametric inference based on undersmoothing or robust bias correction. To simplify the discussion we focus on three concrete cases, which illustrate how the comparisons depend on the available smoothness and kernel order; the messages generalize to any  $S$  and/or  $\kappa$ . For this discussion we let  $\kappa_{\text{us}}$  and  $\kappa_{\text{bc}}$  be the kernel orders used for point estimation in  $I_{\text{us}}$  and  $I_{\text{rbc}}$ , respectively, and restrict attention to sequences  $h \rightarrow 0$  where both confidence intervals are first-order valid, even though robust bias correction allows for a broader bandwidth range. Finally, we set  $\ell = 2$  and  $\bar{\rho} \in (0, \infty)$  based on the above discussion.

For the first case, assume that  $f$  is twice continuously differentiable ( $S = 2$ ) and both methods use second order kernels ( $\kappa_{\text{us}} = \kappa_{\text{bc}} = \ell = 2$ ). In this case, both methods target the *same* bias. The coverage errors for  $I_{\text{us}}$  and  $I_{\text{rbc}}$  then follow directly from Corollaries 1(a) and 2(b) upon plugging in these kernel orders, yielding

$$\left| \mathbb{P}[f \in I_{\text{us}}] - (1 - \alpha) \right| \asymp \frac{1}{nh} + nh^5 + h^2 \quad \text{and} \quad \left| \mathbb{P}[f \in I_{\text{rbc}}] - (1 - \alpha) \right| \asymp \frac{1}{nh} + nh^{5+2\varsigma} + h^{2+\varsigma}.$$

Because  $h \rightarrow 0$  and  $\bar{\rho} \in (0, \infty)$ , the coverage error of  $I_{\text{rbc}}$  vanishes more rapidly by virtue of the bias correction. A higher order kernel ( $\kappa_{\text{us}} > 2$ ) would yield this rate for  $I_{\text{us}}$ .

Second, suppose that the density is four-times continuously differentiable ( $S = 4$ ) but second order kernels are maintained. The relevant results are now Corollaries 1(a) and 2(a). Both methods continue to target the *same* leading bias, but now the additional smoothness



available allows precise characterization of the improvement shown above, and we have

$$|\mathbb{P}[f \in I_{\text{us}}] - (1 - \alpha)| \asymp \frac{1}{nh} + nh^5 + h^2 \quad \text{and} \quad |\mathbb{P}[f \in I_{\text{rbc}}] - (1 - \alpha)| \asymp \frac{1}{nh} + nh^9 + h^4.$$

This case is perhaps the most empirically relevant one, where researchers first choose the order of the kernel (here, second order) and then conduct/optimize inference based on that choice. Indeed, for this case optimal bandwidth choices can be derived (Section 2.4).

Finally, maintain  $S = 4$  but suppose that undersmoothing is based on a fourth-order kernel while bias correction continues to use two second-order kernels ( $k_{\text{us}} = 4$ ,  $k_{\text{bc}} = \ell = 2$ ). This is the exact example given by Hall (1992b, p. 676). Now the two methods target *different* biases, but utilize the *same* amount of smoothness. In this case, the relevant results are again Corollaries 1(a) and 2(a), now with  $k = 4$  and  $k = 2$ , respectively. The two methods have the same coverage error decay rate:

$$|\mathbb{P}[f \in I_{\text{us}}] - (1 - \alpha)| \asymp |\mathbb{P}[f \in I_{\text{rbc}}] - (1 - \alpha)| \asymp \frac{1}{nh} + nh^9 + h^4.$$

Indeed, more can be said: with the notation of Eqn. (4), the difference between  $T_{\text{us}}$  and  $T_{\text{rbc}}$  is the change in “kernel” from  $K$ , a fixed function, to  $M$ , an  $n$ -varying, higher-order kernel, and since  $k_{\text{bc}} + \ell = k_{\text{us}}$ , the two kernels are the same order. [ $M$  acts as a higher-order kernel for bias, but may not strictly fit the definition, as explored in the supplement.] This tight link between undersmoothing and robust bias correction does not carry over straightforwardly to local polynomial regression, as we discussed in more detail in Section 3.

In the context of this final example, it is worth revisiting traditional bias correction. The fact that undersmoothing targets a different, and asymptotically smaller, bias than does explicit bias correction, coupled with the requirement that  $\rho \rightarrow 0$ , implicitly constrains bias correction to remove *less* bias than undersmoothing. This is necessary for traditional bias correction, but on the contrary, robust bias correction attains the *same* coverage error decay rate as undersmoothing under the same assumptions.

In sum, these examples show that under identical assumptions, bias correction is not inferior to undersmoothing and if any additional smoothness is available, can yield improved coverage error. These results are confirmed in our simulations.

## 2.4 Optimal Bandwidth and Data-Driven Choice

The prior sections established that robust bias correction can equal, or outperform, undersmoothing for inference. We now show how the method can be implemented to deliver these results in applications. We mimic typical empirical practice where researchers first choose the order of the kernel, then conduct/optimize inference based on that choice. Therefore, we assume the smoothness is unknown but taken to be large and work within Corollary 2(a), that is, viewing  $k \leq S - 2$  and  $\ell = 2$  as fixed and  $\rho$  bounded and positive. This setup allows characterization of the coverage error optimal bandwidth for robust bias correction.

**Corollary 3.** *Under the conditions of Corollary 2(a) with  $\bar{\rho} \in (0, \infty)$ , if  $h = h_{\text{rbc}}^* = H_{\text{rbc}}^*(\rho)n^{-1/(1+(k+2))}$ , then  $\mathbb{P}[f \in I_{\text{rbc}}] = 1 - \alpha + O(n^{-(k+2)/(1+(k+2))})$ , where*

$$\begin{aligned} H_{\text{rbc}}^*(\bar{\rho}) = \arg \min_{H>0} & \left| H^{-1} q_1(M_{\bar{\rho}}) + H^{1+2(k+2)} (f^{(k+2)})^2 (\mu_{K,k+2} + \bar{\rho}^{-2} \mu_{K,k} \mu_{L,2})^2 q_2(M_{\bar{\rho}}) \right. \\ & \left. + H^{k+2} f^{(k+2)} (\mu_{K,k+2} + \bar{\rho}^{-2} \mu_{K,k} \mu_{L,2}) q_3(M_{\bar{\rho}}) \right|. \end{aligned}$$

We can use this result to give concrete methodological recommendations. At the end of this section we discuss the important issue of interval length. Construction of the interval  $I_{\text{rbc}}$  from Eqn. (5) requires choices of bandwidths  $h$  and  $b$  and kernels  $K$  and  $L$ . Given these choices, the point estimate, bias correction, and variance estimators are then readily computable from data using the formulas above. For the kernels  $K$  and  $L$ , we recommend either second order minimum variance (to minimize interval length) or MSE-optimal kernels (see, e.g., Gasser et al., 1985, and the supplemental appendix).

The bandwidth selections are more important in applications. For the bandwidth  $h$ , Corollary 2(a) shows that the MSE-optimal choice  $h_{\text{mse}}^*$  will deliver valid inference, but will be

suboptimal in general (Corollary 3). From a practical point of view, the robust bias corrected interval  $I_{\text{rbc}}(h)$  is attractive because it allows for the MSE-optimal bandwidth and kernel, and hence is based on the MSE-optimal point estimate, while using the same effective sample for both point estimation and inference. Interestingly, although  $I_{\text{rbc}}(h_{\text{mse}}^*)$  is always valid, its coverage error decays as  $n^{-\min\{4, k+2\}/(1+2k)}$  and is thus rate optimal only for second order kernels ( $k = 2$ ), while otherwise being suboptimal with a slower coverage error rate the larger is the kernel order  $k$ .

Corollary 3 gives the coverage error optimal bandwidth,  $h_{\text{rbc}}^*$ , which can be implemented using a simple direct plug-in (DPI) rule:  $\hat{h}_{\text{dpi}} = \hat{H}_{\text{dpi}} n^{-1/(k+3)}$ , where  $\hat{H}_{\text{dpi}}$  is a plug-in estimate of  $H_{\text{rbc}}^*$  in Corollary 3 formed by replacing the unknown  $f^{(k+2)}$  with a pilot estimate (e.g., a consistent nonparametric estimator based on the appropriate MSE-optimal bandwidth). In the supplement we give precise implementation details, as well as an alternative rule-of-thumb bandwidth selector based on rescaling already available data-driven MSE-optimal choices.

For the bandwidth  $b$ , a simple choice is  $b = h$ , or, equivalently,  $\rho = 1$ . We show in the supplement that setting  $\rho = 1$  has good theoretical properties, minimizing interval length of  $I_{\text{rbc}}$  or the MSE of  $\hat{f}$ , depending on the conditions imposed. In our numerical work, we found that  $\rho = 1$  performed well. As a result, from the practitioner's point of view, the choice of  $b$  (or  $\rho$ ) is completely automatic, leaving only one bandwidth to select.

An extensive simulation study, reported in the supplement, illustrates our findings and explores the numerical performance of these choices. We find that coverage of  $I_{\text{rbc}}$  is robust to both  $h$  and  $\rho$  and that our data-driven bandwidth selectors work well in practice, but we note that estimating bandwidths may have higher-order implications (e.g. Hall and Kang, 2001).

Finally, an important issue in applications is whether the good coverage properties of  $I_{\text{rbc}}$  come at the expense of increased interval length. When coverage is asymptotically correct, Corollaries 1 and 2 show that  $I_{\text{rbc}}$  can accommodate (and will optimally employ) a larger bandwidth (i.e.  $h \rightarrow 0$  more slowly), and hence  $I_{\text{rbc}}$  will have shorter average length in large samples than  $I_{\text{us}}$ . Our simulation study (see below and the supplemental appendix) gives the

same conclusion.

## 2.5 Other Methods of Bias Correction

We study a plug-in bias correction method, but there are alternatives. In particular, as pointed out by a reviewer, a leading alternative is the generalized jackknife method of [Schucany and Sommers \(1977\)](#). We will briefly summarize this approach and show a tight connection to our results (restricting to second-order kernels and  $S \geq 2$  only for simplicity).

The generalized jackknife estimator is  $\hat{f}_{\text{GJ},R} := (\hat{f}_1 - R\hat{f}_2)/(1 - R)$ , where  $\hat{f}_1$  and  $\hat{f}_2$  are two initial kernel density estimators, with possibly different bandwidths  $(h_1, h_2)$  and second-order kernels  $(K_1, K_2)$ . From Eqn. (2), the bias of  $\hat{f}_{\text{GJ},R}$  is  $(1 - R)^{-1}f^{(2)}(h_1^2\mu_{K_1,2} - Rh_2^2\mu_{K_2,2}) + o(h_1^2 + h_2^2)$ , whence choosing  $R = (h_1^2\mu_{K_1,2})/(h_2^2\mu_{K_2,2})$  renders the leading bias term exactly zero. Further, if  $S \geq 4$ ,  $\hat{f}_{\text{GJ},R}$  has bias  $O(h_1^4 + h_2^4)$ ; behaving as a point estimator with  $k = 4$ . To connect this approach to ours, observe that with this choice of  $R$  and  $\tilde{\rho} = h_1/h_2$ , then

$$\hat{f}_{\text{GJ},R} = \frac{1}{nh_1} \sum_{i=1}^n \tilde{M}\left(\frac{X_i - x}{h_1}\right), \quad M(u) = K_1(u) - \tilde{\rho}^{1+2} \left\{ \frac{K_2(\tilde{\rho}u) - \tilde{\rho}^{-1}K_1(u)}{\mu_{K_2,2}(1 - R)} \right\} \mu_{K_1,2},$$

exactly matching Eqn. (4); alternatively, write  $\hat{f}_{\text{GJ},R} = \hat{f}_1 - h_1^2 \tilde{f}^{(2)} \mu_{K_1,2}$ , where

$$\tilde{f}^{(2)} = \frac{1}{nh_2^{1+2}} \sum_{i=1}^n \tilde{L}\left(\frac{X_i - x}{h_2}\right), \quad \tilde{L}(u) = \frac{K_2(u) - \tilde{\rho}^{-1}K_1(\tilde{\rho}^{-1}u)}{\mu_{K_2,2}(1 - R)},$$

is a derivative estimator. Therefore, we can view  $\hat{f}_{\text{GJ},R}$  as a specific kernel  $M$  or a specific derivative estimator, and all our results directly apply to  $\hat{f}_{\text{GJ},R}$ ; hence our paper offers a new way of conducting inference (new Studentization) for this case as well. Though we omit the details to conserve space, this is equally true for local polynomial regression (Section 3).

More generally, our main ideas and generic results apply to many other bias correction methods. For a second example, [Singh \(1977\)](#) also proposed a plug-in bias estimator, but without using the derivative of a kernel. Our results cover this approach as well. See the

supplement for further details and references. The key, common message in all cases is that to improve inference the procedure must account for the additional variability introduced by a bias correction method (i.e., to avoid the mismatch present in  $T_{bc}$ ).

### 3 Local Polynomial Estimation and Inference

This section studies local polynomial regression (Ruppert and Wand, 1994; Fan and Gijbels, 1996), and has two principal aims. First, we show that the conclusions from the density case, and their implications for practice, carry over to odd-degree local polynomials. Second, we show that with proper fixed- $n$  Studentization, coverage error adapts to boundary points. We focus on what is novel relative to the density, chiefly variance estimation and boundary points. For interior points, the implications for coverage error, bandwidth selection, and interval length are all analogous to the density case, and we will not retread those conclusions.

To be specific, throughout this section we focus on the case where the smoothness is large relative to the local polynomial degree  $p$ , which is arguably the most relevant case in practice. The results and discussion in Sections 2.2 and 2.3 carry over, essentially upon changing  $k$  to  $p+1$  and  $\ell$  to  $q-p$  (or  $q-p+1$  for interior points with  $q$  even). Similarly, but with increased notational burden, the conclusions of Section 2.5 also remain true for this section. The present results also extend to multivariate data and derivative estimation.

To begin, we define the regression estimator, its bias, and the bias correction. Given a random sample  $\{(Y_i, X_i) : 1 \leq i \leq n\}$ , the local polynomial estimator of  $m(x) = \mathbb{E}[Y_i | X_i = x]$ , temporarily making explicit the evaluation point, is

$$\hat{m}(x) = \mathbf{e}_0' \hat{\boldsymbol{\beta}}_p, \quad \hat{\boldsymbol{\beta}}_p = \arg \min_{\mathbf{b} \in \mathbb{R}^{p+1}} \sum_{i=1}^n (Y_i - \mathbf{r}_p(X_i - x)' \mathbf{b})^2 K \left( \frac{X_i - x}{h} \right),$$

where, for an integer  $p \geq 1$ ,  $\mathbf{e}_0$  is the  $(p+1)$ -vector with a one in the first position and zeros in the rest, and  $\mathbf{r}_p(u) = (1, u, u^2, \dots, u^p)'$ . We restrict attention to  $p$  odd, as is standard, though the qualifier may be omitted. We define  $\mathbf{Y} = (Y_1, \dots, Y_n)'$ ,  $\mathbf{R}_p = [\mathbf{r}_p((X_1 - x)/h), \dots, \mathbf{r}_p((X_n -$

$x)/h)]'$ ,  $\mathbf{W}_p = \text{diag}(h^{-1}K((X_i - x)/h) : i = 1, \dots, n)$ , and  $\mathbf{\Gamma}_p = \mathbf{R}_p' \mathbf{W}_p \mathbf{R}_p / n$  ( $\text{diag}(a_i : i = 1, \dots, n)$  denotes the  $n \times n$  diagonal matrix constructed using  $a_1, a_2, \dots, a_n$ ). Then, reverting back to omitting the argument  $x$ , the local polynomial estimator is  $\hat{m} = \mathbf{e}_0' \mathbf{\Gamma}_p^{-1} \mathbf{R}_p' \mathbf{W}_p \mathbf{Y} / n$ .

Under regularity conditions below, the conditional bias satisfies

$$\mathbb{E}[\hat{m}|X_1, \dots, X_n] - m = h^{p+1} m^{(p+1)} \frac{1}{(p+1)!} \mathbf{e}_0' \mathbf{\Gamma}_p^{-1} \mathbf{\Lambda}_p + o_P(h^{p+1}), \quad (6)$$

where  $\mathbf{\Lambda}_p = \mathbf{R}_p' \mathbf{W}_p [((X_1 - x)/h)^{p+1}, \dots, ((X_n - x)/h)^{p+1}]' / n$ . Here, the quantity  $\mathbf{e}_0' \mathbf{\Gamma}_p^{-1} \mathbf{\Lambda}_p / (p+1)!$  is random, unlike in the density case (c.f. (2)), but it is known and bounded in probability. Following [Fan and Gijbels \(1996, Section 4.4, p. 116\)](#), we will estimate  $m^{(p+1)}$  in (6) using a second local polynomial regression, of degree  $q > p$  (even or odd), based on a kernel  $L$  and bandwidth  $b$ . Thus,  $\mathbf{r}_q(u)$ ,  $\mathbf{R}_q$ ,  $\mathbf{W}_q$ , and  $\mathbf{\Gamma}_q$  are defined as above, but substituting  $q$ ,  $L$ , and  $b$  in place of  $p$ ,  $K$ , and  $h$ , respectively. Denote by  $\mathbf{e}_{p+1}$  the  $(q+1)$ -vector with one in the  $p+2$  position, and zeros in the rest. Then we estimate the bias with

$$\hat{B}_m = h^{p+1} \hat{m}^{(p+1)} \frac{1}{(p+1)!} \mathbf{e}_0' \mathbf{\Gamma}_p^{-1} \mathbf{\Lambda}_p, \quad \hat{m}^{(p+1)} = b^{-p-1} (p+1)! \mathbf{e}_{p+1}' \mathbf{\Gamma}_q^{-1} \mathbf{R}_q' \mathbf{W}_q \mathbf{Y} / n.$$

Exactly as in the density case,  $\hat{B}_m$  introduces variance that is controlled by  $\rho$  and will be captured by robust bias correction.

### 3.1 Variance Estimation

The Studentizations in the density case were based on fixed- $n$  expectations, and we will show that retaining this is crucial for local polynomials. The fixed- $n$  versus asymptotic distinction is separate from, and more fundamental than, whether we employ feasible versus infeasible quantities. The advantage of fixed- $n$  Studentization also goes beyond bias correction.

To begin, we condition on the covariates so that  $\mathbf{\Gamma}_p^{-1}$  is fixed. Define  $v(\cdot) = \mathbb{V}[Y|X = \cdot]$

and  $\Sigma = \text{diag}(v(X_i) : i = 1, \dots, n)$ . Straightforward calculation gives

$$\sigma_{\text{us}}^2 = (nh)\mathbb{V}[\hat{m}|X_1, \dots, X_n] = \frac{h}{n} \mathbf{e}_0' \Gamma_p^{-1} (\mathbf{R}_p' \mathbf{W}_p \Sigma \mathbf{W}_p \mathbf{R}_p) \Gamma_p^{-1} \mathbf{e}_0. \quad (7)$$

One can then show that  $\sigma_{\text{us}}^2 \rightarrow_P v(x)f(x)^{-1}\mathcal{V}(K, p)$ , with  $\mathcal{V}(K, p)$  a known, constant function of the kernel and polynomial degree. Importantly, both the nonasymptotic form and the convergence hold in the interior or on the boundary, though  $\mathcal{V}(K, p)$  changes.

To first order, one could use  $\sigma_{\text{us}}^2$  or the leading asymptotic term; all that remains is to make each feasible, requiring estimators of the variance function, and for the asymptotic form, also the density. These may be difficult to estimate when  $x$  is a boundary point. Concerned by this, [Chen and Qin \(2002, p. 93\)](#) consider feasible and infeasible versions but conclude that “an increased coverage error near the boundary is still the case even when we know the values of  $f(x)$  and  $v(x)$ .” Our results show that this is not true in general: using fixed- $n$  Studentization, feasible or infeasible, leads to confidence intervals with the same coverage error rates at interior and boundary points, thereby retaining the celebrated boundary carpentry property.

For robust bias correction,  $\sigma_{\text{rbc}}^2 = (nh)V[\hat{m} - \hat{B}_m|X_1, \dots, X_n]$  captures the variance of  $\hat{m}$  and  $\hat{m}^{(p+1)}$  as well as their covariance. A similar fixed- $n$  calculation gives

$$\sigma_{\text{rbc}}^2 = \frac{h}{n} \mathbf{e}_0' \Gamma_p^{-1} (\Xi_{p,q} \Sigma \Xi_{p,q}') \Gamma_p^{-1} \mathbf{e}_0, \quad \Xi_{p,q} = \mathbf{R}_p' \mathbf{W}_p - \rho^{p+2} \Lambda_p \mathbf{e}_{p+1}' \Gamma_q^{-1} \mathbf{R}_q' \mathbf{W}_q \quad (8)$$

To make the fixed- $n$  scalings feasible,  $\hat{\sigma}_{\text{us}}^2$  and  $\hat{\sigma}_{\text{rbc}}^2$  take the forms (7) and (8) and replace  $\Sigma$  with an appropriate estimator. First, we form  $\hat{v}(X_i) = (Y_i - \mathbf{r}_p(X_i - x)' \hat{\beta}_p)^2$  for  $\hat{\sigma}_{\text{us}}^2$  or  $\hat{v}(X_i) = (Y_i - \mathbf{r}_q(X_i - x)' \hat{\beta}_q)^2$  for  $\hat{\sigma}_{\text{rbc}}^2$ . The latter is bias-reduced because  $\mathbf{r}_p(X_i - x)' \beta_p$  is a  $p$ -term Taylor expansion of  $m(X_i)$  around  $x$ , and  $\hat{\beta}_p$  estimates  $\beta_p$  (similarly with  $q$  in place of  $p$ ), and we have  $q > p$ . Next, motivated by the fact that least-squares residuals are on average too small, we appeal to the HCK class of estimators (see [MacKinnon \(2013\)](#) for a review), which are defined as follows. First,  $\hat{\sigma}_{\text{us}}^2\text{-HC0}$  uses  $\hat{\Sigma}_{\text{us}} = \text{diag}(\hat{v}(X_i) : i = 1, \dots, n)$ . Then,  $\hat{\sigma}_{\text{us}}^2\text{-}$

HCk,  $k = 1, 2, 3$ , is obtained by dividing  $\hat{v}(X_i)$  by, respectively,  $(n - 2\text{tr}(\mathbf{Q}_p) + \text{tr}(\mathbf{Q}'_p \mathbf{Q}_p))/n$ ,  $(1 - \mathbf{Q}_{p,ii})$ , or  $(1 - \mathbf{Q}_{p,ii})^2$ , where  $\mathbf{Q}_p := \mathbf{R}'_p \mathbf{\Gamma}_p^{-1} \mathbf{R}'_p \mathbf{W}_p / n$  is the projection matrix and  $\mathbf{Q}_{p,ii}$  its  $i$ -th diagonal element. The corresponding estimators  $\hat{\sigma}_{\text{rbc}}^2$ -HCk are the same, but with  $q$  in place of  $p$ . For theoretical results, we use HC0 for concreteness and simplicity, though inspection of the proof shows that simple modifications allow for the other HCk estimators and rates do not change. These estimators may perform better for small sample sizes. Another option is to use a nearest-neighbor-based variance estimators with a fixed number of neighbors, following the ideas of [Muller and Stadtmuller \(1987\)](#) and [Abadie and Imbens \(2008\)](#). Note that none of these estimators assume local or global homoskedasticity nor rely on new tuning parameters. Details and simulation results for all these estimators are given in the supplement, see §S.II.2.3 and Table S.II.9.

### 3.2 Higher Order Expansions of Coverage Error

Recycling notation to emphasize the parallel, we study the following three statistics:

$$T_{\text{us}} = \frac{\sqrt{nh}(\hat{m} - m)}{\hat{\sigma}_{\text{us}}}, \quad T_{\text{bc}} = \frac{\sqrt{nh}(\hat{m} - \hat{B}_m - m)}{\hat{\sigma}_{\text{us}}}, \quad T_{\text{rbc}} = \frac{\sqrt{nh}(\hat{m} - \hat{B}_m - m)}{\hat{\sigma}_{\text{rbc}}},$$

and their associated confidence intervals  $I_{\text{us}}$ ,  $I_{\text{bc}}$ , and  $I_{\text{rbc}}$ , exactly as in Eqn. (5). Importantly, all present definitions and results are valid for an evaluation point in the interior and at the boundary of the support of  $X_i$ . The following standard conditions will suffice, augmented with the appropriate Cramér's condition given in the supplement to conserve space.

**Assumption 3.1** (Data-generating process).  *$\{(Y_1, X_1), \dots, (Y_n, X_n)\}$  is a random sample, where  $X_i$  has the absolutely continuous distribution with Lebesgue density  $f$ ,  $\mathbb{E}[Y^{8+\delta}|X] < \infty$  for some  $\delta > 0$ , and in a neighborhood of  $x$ ,  $f$  and  $v$  are continuous and bounded away from zero,  $m$  is  $S > q + 2$  times continuously differentiable with bounded derivatives, and  $m^{(S)}$  is Hölder continuous with exponent  $\varsigma$ .*



**Assumption 3.2** (Kernels). *The kernels  $K$  and  $L$  are positive, bounded, even functions, and with compact support.*

We now give our main, generic result for local polynomials, analogous to Theorem 1. For notation, the polynomials  $q_1$ ,  $q_2$ , and  $q_3$  and the biases  $\eta_{\text{us}}$  and  $\eta_{\text{bc}}$ , are cumbersome and exact forms are deferred to the supplement. All that matters is that the polynomials are known, odd, bounded, and bounded away from zero and that the biases have the usual convergence rates, as detailed below.

**Theorem 2.** *Let Assumptions 3.1, 3.2, and Cramér's condition hold and  $nh/\log(nh) \rightarrow \infty$ .*

(a) *If  $\eta_{\text{us}} \log(nh) \rightarrow 0$ , then*

$$\mathbb{P}[m \in I_{\text{us}}] = 1 - \alpha + \left\{ \frac{1}{nh} q_{1,\text{us}} + \eta_{\text{us}}^2 q_{2,\text{us}} + \frac{\eta_{\text{us}}}{\sqrt{nh}} q_{3,\text{us}} \right\} \phi(z_{\frac{\alpha}{2}}) \{1 + o(1)\}.$$

(b) *If  $\eta_{\text{bc}} \log(nh) \rightarrow 0$  and  $\rho \rightarrow 0$ , then*

$$\begin{aligned} \mathbb{P}[m \in I_{\text{bc}}] = 1 - \alpha + \left\{ \frac{1}{nh} q_{1,\text{us}} + \eta_{\text{bc}}^2 q_{2,\text{us}} + \frac{\eta_{\text{bc}}}{\sqrt{nh}} q_{3,\text{us}} \right\} \phi(z_{\frac{\alpha}{2}}) \{1 + o(1)\} \\ + \rho^{p+2} (\Omega_{1,\text{bc}} + \rho^{p+1} \Omega_{2,\text{bc}}) \phi(z_{\frac{\alpha}{2}}) z_{\frac{\alpha}{2}} \{1 + o(1)\}. \end{aligned}$$

(c) *If  $\eta_{\text{bc}} \log(nh) \rightarrow 0$  and  $\rho \rightarrow \bar{\rho} < \infty$ , then*

$$\mathbb{P}[m \in I_{\text{rbc}}] = 1 - \alpha + \left\{ \frac{1}{nh} q_{1,\text{rbc}} + \eta_{\text{bc}}^2 q_{2,\text{rbc}} + \frac{\eta_{\text{bc}}}{\sqrt{nh}} q_{3,\text{rbc}} \right\} \phi(z_{\frac{\alpha}{2}}) \{1 + o(1)\}.$$

This theorem, which covers both interior and boundary points, establishes that the conclusions found in the density case carry over to odd-degree local polynomial regression. (Although we focus on  $p$  odd, part (a) is valid in general and (b) and (c) are valid at the boundary for  $p$  even.) In particular, this shows that robust bias correction is as good as, or better than, undersmoothing in terms of coverage error. Traditional bias correction is again inferior due to the variance and covariance terms  $\rho^{p+2}(\Omega_{1,\text{bc}} + \rho^{p+1}\Omega_{2,\text{bc}})$ . Coverage error optimal bandwidths can be derived as well, and similar conclusions are found. Best possible rates are defined for fixed  $p$  here, the analogue of  $k$  above; see Section 2.2 for further discussion on smoothness.

Before discussing bias correction, one aspect of the undersmoothing result is worth mentioning. The fact that Theorem 2 covers both interior and boundary points, without requiring additional assumptions, is in some sense, expected: one of the strengths of local polynomial estimation is its adaptability to boundary points. In particular, from Eqn. (6) and  $p$  odd it follows that  $\eta_{\text{us}} \asymp \sqrt{nh}h^{p+1}$  at the interior and the boundary. Therefore, part (a) shows that the decay rate in coverage error does not change at the boundary for the standard confidence interval (but the leading constants will change). This finding contrasts with the result of Chen and Qin (2002) who studied the special case  $p = 1$  without bias correction (part (a) of Theorem 2), and is due entirely to the fixed- $n$  Studentization.

Turning to robust bias correction, we will, in contrast, find rate differences between the interior and the boundary, no matter the parity of  $q$ . As before,  $\eta_{\text{bc}}$  has two terms, representing the higher-order bias of the point estimator and the bias of the bias estimator. The former can be viewed as the bias if  $m^{(p+1)}$  were zero, and since  $p + 1$  is even, we find that it is of order  $\sqrt{nh}h^{p+3}$  in the interior but  $\sqrt{nh}h^{p+2}$  at the boundary. The bias of the bias correction depends on both bandwidths  $h$  and  $b$ , as well as  $p$  and  $q$ , in exact analogy to the density case. For  $q$  odd, it is of order  $h^{p+1}b^{q-p}$  at all points, whereas for  $q$  even this rate is attained at the boundary, but in the interior the order increases to  $h^{p+1}b^{q+1-p}$ . Collecting these facts: in the interior,  $\eta_{\text{bc}} \asymp \sqrt{nh}h^{p+3}(1 + \rho^{-2}b^{q-p-2})$  for odd  $q$  or with  $b^{q-p-1}$  for  $q$  even; at the boundary,  $\eta_{\text{bc}} \asymp \sqrt{nh}h^{p+2}(1 + \rho^{-1}b^{q-p-1})$ . Further details are in the supplement.

In light of these rates, the same logic of Section 2.2 leads us to restrict attention to bounded, positive  $\rho$  and  $q = p + 1$ , and thus even. Calonico et al. (2014, Remark 7) point out that in the special case of  $q = p + 1$ ,  $K = L$ , and  $\rho = 1$ ,  $\hat{m} - \hat{B}_m$  is identical to a local polynomial estimator of order  $q$ ; this is the closest analogue to  $M$  being a higher-order kernel. If the point of interest is in the interior, then  $q = p + 2$  yields the same rates.

For notational ease, let  $\tilde{\eta}_{\text{bc}}^{\text{int}}$  and  $\tilde{\eta}_{\text{bc}}^{\text{bnd}}$  be the leading constants for the interior and boundary, respectively, so that e.g.  $\eta_{\text{bc}} = \sqrt{nh}h^{p+3}[\tilde{\eta}_{\text{bc}}^{\text{int}} + o(1)]$  in the interior (exact expressions are in the supplement). We then have the following, precise result; the analogue of Corollary 2(a).

**Corollary 4.** *Let the conditions of Theorem 2(c) hold, with  $\bar{\rho} \in (0, \infty)$  and  $q = p + 1$ .*

(a) *For an interior point,*

$$\mathbb{P}[m \in I_{\text{rbc}}] = 1 - \alpha + \left\{ \frac{1}{nh} q_{1,\text{rbc}} + nh^{1+2(p+3)} (\tilde{\eta}_{\text{bc}}^{\text{int}})^2 q_{2,\text{rbc}} + h^{p+3} (\tilde{\eta}_{\text{bc}}^{\text{int}}) q_{3,\text{rbc}} \right\} \phi(z_{\frac{\alpha}{2}}) \{1 + o(1)\}.$$

(b) *For a boundary point,*

$$\mathbb{P}[m \in I_{\text{rbc}}] = 1 - \alpha + \left\{ \frac{1}{nh} q_{1,\text{rbc}} + nh^{1+2(p+2)} (\tilde{\eta}_{\text{bc}}^{\text{bnd}})^2 q_{2,\text{rbc}} + h^{p+2} (\tilde{\eta}_{\text{bc}}^{\text{bnd}}) q_{3,\text{rbc}} \right\} \phi(z_{\frac{\alpha}{2}}) \{1 + o(1)\}.$$

There are differences in both the rates and constants between parts (a) and (b) of this result, though most of the changes to constants are “hidden” notationally by the definitions of  $\tilde{\eta}_{\text{bc}}^{\text{bnd}}$  and the polynomials  $q_{k,\text{rbc}}$ . Part (a) most closely resembles Corollary 2 due to the symmetry yielding the corresponding rate improvement (recall that  $k$  in the density case is replaced with  $p + 1$  here), and hence all the corresponding conclusions hold qualitatively for local polynomials.

### 3.3 Practical Choices and Empirical Consequences

As we did for the density, we now derive bandwidth choices, and data-driven implementations, to optimize coverage error in applications.

**Corollary 5.** *Let the conditions of Corollary 4 hold.*

(a) *For an interior point, if  $h = h_{\text{rbc}}^* = H_{\text{rbc}}^* n^{-1/(p+4)}$ , then  $\mathbb{P}[m \in I_{\text{rbc}}] = 1 - \alpha + O(n^{-(p+3)/(p+4)})$ , where*

$$H_{\text{rbc}}^* = \arg \min_{H>0} |H^{-1} q_{1,\text{rbc}} + H^{1+2(p+3)} (\tilde{\eta}_{\text{bc}}^{\text{int}})^2 q_{2,\text{rbc}} + H^{p+3} (\tilde{\eta}_{\text{bc}}^{\text{int}}) q_{3,\text{rbc}}|.$$

(b) *For a boundary point, if  $h = h_{\text{rbc}}^* = H_{\text{rbc}}^*(\rho) n^{-1/(p+3)}$ , then  $\mathbb{P}[m \in I_{\text{rbc}}] = 1 - \alpha + O(n^{-(p+2)/(p+3)})$ , where*

$$H_{\text{rbc}}^*(\bar{\rho}) = \arg \min_{H>0} |H^{-1} q_{1,\text{rbc}} + H^{1+2(p+2)} (\tilde{\eta}_{\text{bc}}^{\text{bnd}})^2 q_{2,\text{rbc}} + H^{p+2} (\tilde{\eta}_{\text{bc}}^{\text{bnd}}) q_{3,\text{rbc}}|$$

To implement these results, we first set  $\rho = 1$  and the kernels  $K$  and  $L$  equal to any desired second order kernel, typical choices being triangular, Epanechnikov, and uniform. The variance estimator  $\hat{\sigma}_{\text{rbc}}$  is defined in Section 3.1, and is fully implementable, and thus so is  $I_{\text{rbc}}$ , once the bandwidth  $h$  is chosen.

For selecting  $h$  at an interior point, the same conclusions from density estimation apply: (i) coverage of  $I_{\text{rbc}}$  is quite robust with respect to  $h$  and  $\rho$ , (ii) feasible choices for  $h$  are easy to construct, and (iii) an MSE-optimal bandwidth only delivers the best coverage error for  $p = 1$  (that is,  $k = 2$  in the density case). On the other hand, for a boundary point, an interesting consequence of Corollary 5 is that an MSE-optimal bandwidth *never* delivers optimal coverage error decay rates, even for local linear regression:  $h_{\text{mse}}^* \propto n^{-1/(2p+3)} \gg h_{\text{rbc}}^* \propto n^{-1/(p+3)}$ .

Keeping this in mind, we give a fully data-driven direct plug-in (DPI) bandwidth selector for both interior and boundary points:  $\hat{h}_{\text{dpi}}^{\text{int}} = \hat{H}_{\text{dpi}}^{\text{int}} n^{-1/(p+4)}$  and  $\hat{h}_{\text{dpi}}^{\text{bnd}} = \hat{H}_{\text{dpi}}^{\text{bnd}} n^{-1/(p+3)}$ , where  $\hat{H}_{\text{dpi}}^{\text{int}}$  and  $\hat{H}_{\text{dpi}}^{\text{bnd}}$  are estimates of (the appropriate)  $H_{\text{rbc}}^*$  of Corollary 5, obtained by estimating unknowns by pilot estimators employing a readily-available pilot bandwidth. The complete steps to form  $\hat{H}_{\text{dpi}}^{\text{int}}$  and  $\hat{H}_{\text{dpi}}^{\text{bnd}}$  are in the supplement, as is a second data-driven bandwidth choice, based on rescaling already-available MSE-optimal bandwidths. All our methods are available in the R package `nprobust`; see <https://cran.r-project.org/package=nprobust>.

## 4 Simulation Results

We now report a representative sample of results from a simulation study to illustrate our findings. We drew 5,000 replicated data sets, each being  $n = 500$  i.i.d. draws from the model  $Y_i = m(X_i) + \varepsilon_i$ , with  $m(x) = \sin(3\pi x/2)(1 + 18x^2[\text{sgn}(x) + 1])^{-1}$ ,  $X_i \sim \mathcal{U}[0, 1]$ , and  $\varepsilon_i \sim \mathcal{N}(0, 1)$ . We consider inference at the five points  $x \in \{-2/3, -1/3, 0, 1/3, 2/3\}$ . The function  $m(x)$  and the five evaluation points are plotted in Figure 1; this function was previously used by Berry et al. (2002) and Hall and Horowitz (2013). The supplement gives results for other models, bandwidth selectors and their simulation distributions, alternative

variance estimators, and more detailed studies of coverage and length.

We compared robust bias correction to undersmoothing, traditional bias correction, the off-the-shelf R package `locfit` (Loader, 2013), and the procedure of Hall and Horowitz (2013). In all cases the point estimator is based on local linear regression with the data-driven bandwidth  $\hat{h}_{\text{dpi}}^{\text{int}}$ , which shares the rate of  $\hat{h}_{\text{mse}}$  in this case, and  $\rho = 1$ . The `locfit` package has a bandwidth selector, but it was ill-behaved and often gave zero empirical coverage. Hall and Horowitz (2013) do not give an explicit optimal bandwidth, but do advocate a feasible  $\hat{h}_{\text{mse}}$ , following Ruppert et al. (1995). To implement their method, we used 500 bootstrap replications and we set  $1 - \xi = 0.9$  over a sequence  $\{x_1, \dots, x_N\} = \{-0.9, -0.8, \dots, 0, \dots, 0.8, 0.9\}$  to obtain the final quantile  $\hat{\alpha}_\xi(\alpha_0)$ , and used their proposed standard errors  $\hat{\sigma}_{\text{HH}}^2 = \kappa \hat{\sigma}^2 / \hat{f}_X$ , where  $\hat{\sigma}^2 = \sum_{i=1}^n \hat{\varepsilon}_i^2 / n$  for  $\hat{\varepsilon}_i = \tilde{\varepsilon}_i - \bar{\varepsilon}$ , with  $\tilde{\varepsilon}_i = Y_i - \hat{m}(X_i)$  and  $\bar{\varepsilon} = \sum_{i=1}^n \tilde{\varepsilon}_i / n$ .

Table 1 shows empirical coverage and average length at all five points for all five methods. Robust bias correction yields accurate coverage throughout the support; performance of the other methods varies. For  $x = -2/3$ , the regression function is nearly linear, leaving almost no bias, and the other methods work quite well. In contrast, at  $x = -1/3$  and  $x = 0$ , all methods except robust bias correction suffer from coverage distortions due to bias. Indeed, Hall and Horowitz (2013, p. 1893) report that “[t]he ‘exceptional’ 100% of points that are not covered are typically close to the locations of peaks and troughs, [which] cause difficulties because of bias.” Finally, bias is still present, though less of a problem, for  $x = 1/3$  and  $x = 2/3$ , and coverage of the competing procedures improves somewhat. Motivated by the fact that the data-driven bandwidth selectors may be “too large” for proper undersmoothing, we studied the common practice of ad-hoc undersmoothing of the MSE-optimal bandwidth choice  $\hat{h}_{\text{mse}}$ : the results in Table S.II.8 of the supplement show this to be no panacea.

To illustrate our findings further, Figures 2(a) and 2(b) compare coverage and length of different inference methods over a range of bandwidths. Robust bias correction delivers accurate coverage for a wide range of bandwidths, including larger choices, and thus can yield shorter intervals. For undersmoothing, coverage accuracy requires a delicate choice of

bandwidth, and for correct coverage, a longer interval. Figure 2(c), in color online, reinforces this point by showing the “average position” of  $I_{\text{us}}(h)$  and  $I_{\text{rbc}}(h)$  for a range of bandwidths: each bar is centered at the average bias and is of average length, and then color-coded by coverage (green indicates good coverage, fading to red as coverage deteriorates). These results show that when  $I_{\text{us}}$  is short, bias is large and coverage is poor. In contrast,  $I_{\text{rbc}}$  has good coverage at larger bandwidths and thus shorter length.

## 5 Conclusion

This paper has made three distinct, but related points regarding nonparametric inference. First, we showed that bias correction, when coupled with a new standard error formula, performs as well or better than undersmoothing for confidence interval coverage and length. Further, such intervals are more robust to bandwidth choice in applications. Second, we showed theoretically when the popular empirical practice of using MSE-optimal bandwidths is justified, and more importantly, when it is not, and we gave concrete implementation recommendations for applications. Third, we proved that confidence intervals based on local polynomials do have automatic boundary carpentry, provided proper Studentization is used. These results are tied together through the themes of higher order expansions and the importance of finite sample variance calculations and the key, common message that inference procedures must account for additional variability introduced by bias correction.

## 6 References

- Abadie, A., and Imbens, G. W. (2008), “Estimation of the Conditional Variance in Paired Experiments,” *Annales d’Economie et de Statistique*, 175–187.
- Andrews, D. W. K. (2002), “Higher-Order Improvements of a Computationally Attractive  $k$ -Step Bootstrap for Extremum Estimators,” *Econometrica*, 70, 119–162.
- Armstrong, T. B., and Kolesár, M. (2015), “Optimal inference in a class of regression models,” *Arxiv preprint arXiv:1511.06028*.

- Berry, S. M., Carroll, R. J., and Ruppert, D. (2002), “Bayesian Smoothing and Regression Splines for Measurement Error Problems,” *Journal of the American Statistical Association*, 97, 160–169.
- Bhattacharya, R. N., and Ghosh, J. K. (1978), “On the Validity of the Formal Edgeworth Expansion,” *The Annals of Statistics*, 6, 434–451.
- Bhattacharya, R. N., and Rao, R. R. (1976), *Normal Approximation and Asymptotic Expansions*, John Wiley and Sons.
- Calonico, S., Cattaneo, M. D., and Farrell, M. H. (2016), “Coverage Error Optimal Confidence Intervals for Regression Discontinuity Designs,” *working paper*.
- Calonico, S., Cattaneo, M. D., and Titiunik, R. (2014), “Robust Nonparametric Confidence Intervals for Regression-Discontinuity Designs,” *Econometrica*, 82, 2295–2326.
- Cattaneo, M. D., and Farrell, M. H. (2013), “Optimal Convergence Rates, Bahadur Representation, and Asymptotic Normality of Partitioning Estimators,” *Journal of Econometrics*, 174, 127–143.
- Chen, S. X., and Qin, Y. S. (2002), “Confidence Intervals Based on Local Linear Smoother,” *Scandinavian Journal of Statistics*, 29, 89–99.
- Fan, J., and Gijbels, I. (1996), *Local polynomial modelling and its applications*, London: Chapman and Hall.
- Gasser, T., Muller, H.-G., and Mammitzsch, V. (1985), “Kernels for Nonparametric Curve Estimation,” *Journal of the Royal Statistical Society. Series B*, 47, 238–252.
- Hall, P. (1991), “Edgeworth Expansions for Nonparametric Density Estimators, with Applications,” *Statistics*, 22, 215–232.
- (1992a), *The Bootstrap and Edgeworth Expansion*, New York: Springer-Verlag.
- (1992b), “Effect of Bias Estimation on Coverage Accuracy of Bootstrap Confidence Intervals for a Probability Density,” *The Annals of Statistics*, 20, 675–694.
- Hall, P., and Horowitz, J. L. (2013), “A Simple Bootstrap Method for Constructing Nonparametric Confidence Bands for Functions,” *The Annals of Statistics*, 41, 1892–1921.
- Hall, P., and Kang, K.-H. (2001), “Bootstrapping Nonparametric Density Estimators with Empirically Chosen Bandwidths,” *The Annals of Statistics*, 29, 1443–1468.
- Hansen, B. E. (2015), “Robust Inference,” *working paper*, *University of Wisconsin*.
- Horowitz, J. L. (2009), *Semiparametric and Nonparametric Methods in Econometrics*, Springer.
- James, G. S., and Mayne, A. J. (1962), “Cumulants of Functions of Random Variables,” *Sankhyā*, 24, 47–54.

- Jones, M. C. (1994), “On Kernel Density Derivative Estimation,” *Communications in Statistics - Theory and Methods*, 23, 2133–2139.
- Jones, M. C., and Foster, P. J. (1993), “Generalized Jackknifing and Higher Order Kernels,” *Journal of Nonparametric Statistics*, 3, 81–94.
- Jones, M. C., and Signorini, D. F. (1997), “A Comparison of Higher-Order Bias Kernel Density Estimators,” *Journal of the American Statistical Association*, 92, 1063–1073.
- Kline, P., and Santos, A. (2012), “Higher order properties of the wild bootstrap under misspecification,” *Journal of Econometrics*, 171, 54–70.
- Loader, C. (2013), *locfit: Local Regression, Likelihood and Density Estimation.*, R package version 1.5-9.1.
- MacKinnon, J. G. (2013), “Thirty Years of Heteroskedasticity-Robust Inference,” in *Recent Advances and Future Directions in Causality, Prediction, and Specification Analysis*, eds. X. Chen and N. R. Swanson, Springer, pp. 437–461.
- Marron, J. S., and Wand, M. P. (1992), “Exact Mean Integrated Squared Error,” *The Annals of Statistics*, 20, 712–736.
- Muller, H.-G., and Stadtmüller, U. (1987), “Estimation of Heteroscedasticity in Regression Analysis,” *The Annals of Statistics*, 15, 610–625.
- Mykland, P., and Zhang, L. (2015), “Assessment of Uncertainty in High Frequency Data: The Observed Asymptotic Variance,” *Econometrica*, *forthcoming*.
- Neumann, M. H. (1997), “Pointwise confidence intervals in nonparametric regression with heteroscedastic error structure,” *Statistics*, 29, 1–36.
- Ruppert, D., Sheather, S. J., and Wand, M. P. (1995), “An Effective Bandwidth Selector for Local Least Squares Regression,” *Journal of the American Statistical Association*, 90, 1257–1270.
- Ruppert, D., and Wand, M. P. (1994), “Multivariate Locally Weighted Least Squares Regression,” *The Annals of Statistics*, 22, 1346–1370.
- Ruppert, D., Wand, M. P., and Carroll, R. (2009), *Semiparametric Regression*, New York: Cambridge University Press.
- Schennach, S. M. (2015), “A bias bound approach to nonparametric inference,” *cemmap working paper CWP71/15*.
- Schucany, W., and Sommers, J. P. (1977), “Improvement of Kernel Type Density Estimators,” *Journal of the American Statistical Association*, 72, 420–423.
- Silverman, B. W. (1986), *Density Estimation for Statistics and Data Analysis*, Chapman & Hall.



- Singh, R. S. (1977), “Improvement on Some Known Nonparametric Uniformly Consistent Estimators of Derivatives of a Density,” *The Annals of Statistics*, 5, 394–399.
- Stone, C. J. (1982), “Optimal Global Rates of Convergence for Nonparametric Regression,” *The Annals of Statistics*, 10, 1040–1053.
- Wand, M., and Jones, M. (1995), *Kernel Smoothing*, Florida: Chapman & Hall/CRC.

Table 1: Empirical Coverage and Average Interval Length of 95% Confidence Intervals

Evaluation Point	Average Bandwidth	Empirical Coverage					Interval Length			
		US	Locfit	BC	HH	RBC	US	Locfit	HH	RBC
-2/3	0.166	94.8	94.4	81.8	93.5	93.7	0.505	0.544	0.479	0.722
-1/3	0.283	56.5	70.7	80.6	48.2	92.8	0.380	0.409	0.316	0.540
0	0.318	74.4	83.7	80.3	61.1	92.6	0.354	0.383	0.279	0.507
1/3	0.370	89.9	92.1	78.5	78.4	92.9	0.327	0.356	0.241	0.470
2/3	0.265	93.9	93.9	81.3	88.4	93.6	0.391	0.425	0.339	0.562

**Notes:** (i) Column “Average Bandwidth” reports simulation average of estimated bandwidths  $h = \hat{h}_{\text{dpi}} \equiv \hat{h}_{\text{dpi}}^{\text{int}}$ . Simulation distributions for estimated bandwidths are reported in the supplement. (ii) US = Under-smoothing, Locfit = R package `locfit` by Loader (2013), BC = Bias Corrected, HH = Hall and Horowitz (2013), RBC = Robust Bias Corrected.

Figure 1: True Regression Model and Evaluation Points

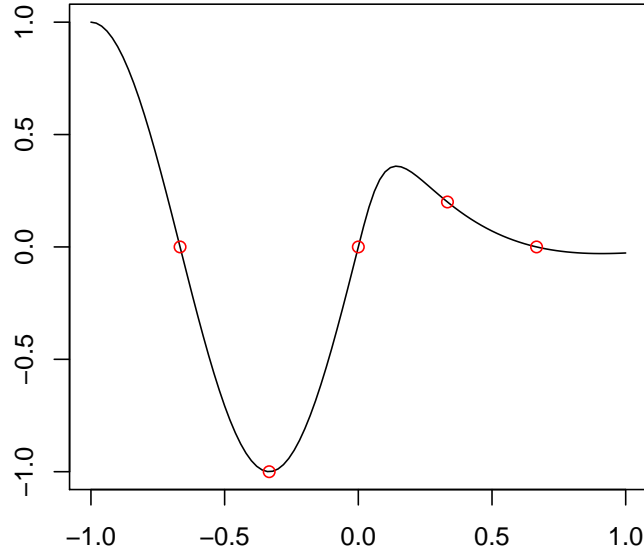
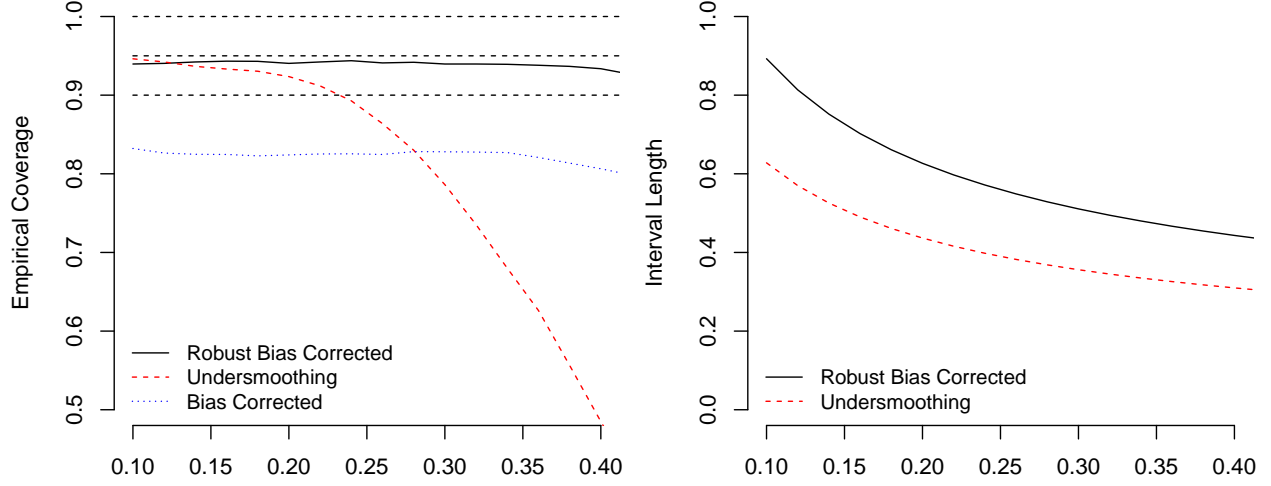
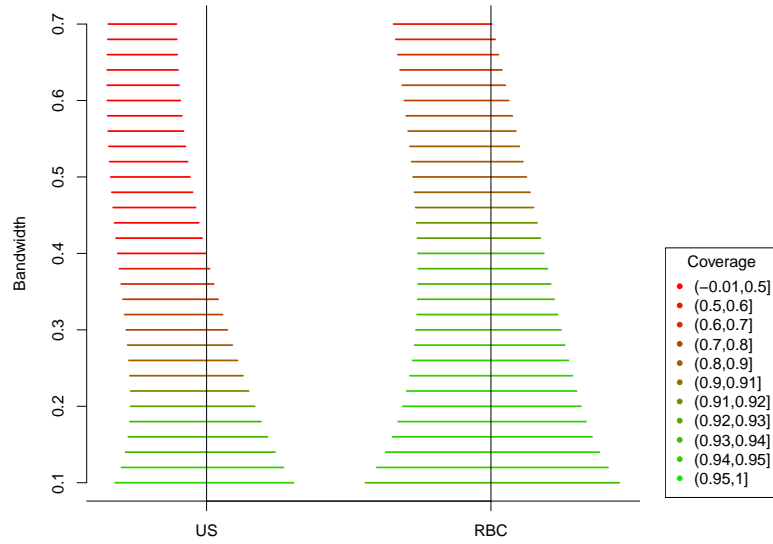


Figure 2: Local Polynomial Simulation Results for  $x = 0$



(a) Empirical Coverage

(b) Interval Length



(c) Average Confidence Intervals and Coverage for Undersmoothing (US) and Robust Bias Correction (RBC).

# Supplement to “On the Effect of Bias Estimation on Coverage Accuracy in Nonparametric Inference”

This supplement contains technical and notational details omitted from the main text, proofs of all results, further technical details and derivations, and additional simulations results and numerical analyses. The main results are Edgeworth expansions of the distribution functions of the  $t$ -statistics  $T_{\text{us}}$ ,  $T_{\text{bc}}$ , and  $T_{\text{rbc}}$ , for density estimation and local polynomial regression. Stating and proving these results is the central purpose of this supplement. The higher-order expansions of confidence interval coverage probabilities in the main paper follow immediately by evaluating the Edgeworth expansions at the interval endpoints.

Part [S.I](#) contains all material for density estimation at interior points, while Part [S.II](#) treats local polynomial regression at both interior and boundary points, as in the main text. Roughly, these have the same generic outline:

- We first present all notation, both for the estimators themselves and the Edgeworth expansions, regardless of when the notation is used, as a collective reference;
- We then discuss optimal bandwidths and other practical matters, expanding on details of the main text;
- Assumptions for validity of the Edgeworth expansions are restated from the main text, and Cramér’s condition is discussed;
- Bias properties are discussed in more detail than in the main text, and some things mentioned there are made precise;
- The main Edgeworth expansions are stated, some corollaries are given, and the proofs are given;
- Complete simulation results are presented.

All our methods are implemented in software available from the authors’ websites and via the R package `nprobust` available at <https://cran.r-project.org/package=nprobust>.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Density Estimation and Inference</b>	<b>5</b>
2.1	Generic Higher Order Expansions of Coverage Error . . . . .	8
2.2	Coverage Error and the Role of Smoothness . . . . .	11
2.3	Comparing Undersmoothing and Robust Bias Correction . . . . .	14
2.4	Optimal Bandwidth and Data-Driven Choice . . . . .	16
2.5	Other Methods of Bias Correction . . . . .	18
<b>3</b>	<b>Local Polynomial Estimation and Inference</b>	<b>19</b>
3.1	Variance Estimation . . . . .	20
3.2	Higher Order Expansions of Coverage Error . . . . .	22
3.3	Practical Choices and Empirical Consequences . . . . .	25
<b>4</b>	<b>Simulation Results</b>	<b>26</b>
<b>5</b>	<b>Conclusion</b>	<b>28</b>
<b>6</b>	<b>References</b>	<b>28</b>
<b>S.I</b>	<b>Kernel Density Estimation and Inference</b>	<b>3</b>
<b>S.I.1</b>	<b>Notation</b>	<b>3</b>
S.I.1.1	Estimators, Variances, and Studentized Statistics . . . . .	3
S.I.1.2	Edgeworth Expansion Terms . . . . .	4
<b>S.I.2</b>	<b>Details of practical implementation</b>	<b>5</b>
S.I.2.1	Bandwidth Choice: Rule-of-Thumb (ROT) . . . . .	6
S.I.2.2	Bandwidth Choice: Direct Plug-In (DPI) . . . . .	6
S.I.2.3	Choice of $\rho$ . . . . .	8
<b>S.I.3</b>	<b>Assumptions</b>	<b>10</b>
<b>S.I.4</b>	<b>Bias</b>	<b>11</b>
S.I.4.1	Precise Bias Calculations . . . . .	11
S.I.4.2	Properties of the kernel $M_\rho(\cdot)$ . . . . .	13
S.I.4.3	Other Bias Reduction Methods . . . . .	15
<b>S.I.5</b>	<b>First Order Properties</b>	<b>17</b>
<b>S.I.6</b>	<b>Main Result: Edgeworth Expansion</b>	<b>18</b>
S.I.6.1	Undersmoothing vs. Bias-Correction Exhausting all Smoothness . . . .	20
S.I.6.2	Multivariate Densities and Derivative Estimation . . . . .	21
<b>S.I.7</b>	<b>Proof of Main Result</b>	<b>24</b>
S.I.7.1	Computing the Terms of the Expansion . . . . .	27
<b>S.I.8</b>	<b>Complete Simulation Results</b>	<b>30</b>
<b>S.II</b>	<b>Local Polynomial Estimation and Inference</b>	<b>47</b>
<b>S.II.1</b>	<b>Notation</b>	<b>47</b>
S.II.1.1	Estimators, Variances, and Studentized Statistics . . . . .	48
S.II.1.2	Edgeworth Expansion Terms . . . . .	50

<b>S.II.2 Details of Practical Implementation</b>	<b>55</b>
S.II.2.1 Bandwidth Choice: Rule-of-Thumb (ROT) . . . . .	55
S.II.2.2 Bandwidth Choice: Direct Plug-In (DPI) . . . . .	55
S.II.2.3 Alternative Standard Errors . . . . .	58
<b>S.II.3 Assumptions</b>	<b>59</b>
<b>S.II.4 Bias</b>	<b>62</b>
<b>S.II.5 Main Result: Edgeworth Expansion</b>	<b>64</b>
S.II.5.1 Coverage Error for Undersmoothing . . . . .	65
<b>S.II.6 Proof of Main Result</b>	<b>66</b>
S.II.6.1 Proof of Theorem 5(a) . . . . .	66
S.II.6.2 Proof of Theorem 5(b) & (c) . . . . .	71
S.II.6.3 Lemmas . . . . .	73
S.II.6.4 Computing the Terms of the Expansion . . . . .	83
<b>S.II.7 Complete Simulation Results</b>	<b>88</b>

## Part S.I

# Kernel Density Estimation and Inference

### S.I.1 Notation

Here we collect notation to be used throughout this section, even if it is restated later. Throughout this supplement, let  $X_{h,i} = (x - X_i)/h$  and similarly for  $X_{b,i}$ . The evaluation point is implicit here. In the course of proofs we will frequently write  $s = \sqrt{nh}$ .

#### S.I.1.1 Estimators, Variances, and Studentized Statistics

To begin, recall that the original and bias-corrected density estimators are

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K(X_{h,i})$$

and

$$\hat{f} - \hat{B}_f = \frac{1}{nh} \sum_{i=1}^n M(X_{h,i}), \quad M(u) := K(u) - \rho^{1+k} L^{(k)}(\rho u) \mu_{K,k}, \quad (9)$$

for symmetric kernel functions  $K(\cdot)$  and  $L(\cdot)$  that integrate to one on their compact support,  $h$  and  $b$  are bandwidth sequences that vanish as  $n \rightarrow \infty$ , and where

$$\rho = h/b, \quad \hat{B}_f = h^k \hat{f}^{(k)}(x) \mu_{K,k}, \quad \hat{f}^{(k)}(x) = \frac{1}{nb^{1+k}} \sum_{i=1}^n L^{(k)}(X_{b,i}),$$

and integrals of the kernel are denoted

$$\mu_{K,k} = \frac{(-1)^k}{k!} \int u^k K(u) du, \quad \text{and} \quad \vartheta_{K,k} = \int K(u)^k du.$$

The three statistics  $T_{\text{us}}$ ,  $T_{\text{bc}}$ , and  $T_{\text{rbc}}$  share a common structure that is exploited to give a unified theorem statement and proof. For  $v \in \{1, 2\}$ , define

$$\hat{f}_v = \frac{1}{nh} \sum_{i=1}^n N_v(X_{h,i}), \quad \text{where} \quad N_1(u) = K(u) \text{ and } N_2(u) = M(u),$$

and  $M$  is given in Eqn. (9). Thus,  $\hat{f}_1 = \hat{f}$  and  $\hat{f}_2 = \hat{f} - \hat{B}_f$ . In exactly the same way, define

$$\sigma_v^2 := nh\mathbb{V}[\hat{f}_v] = \frac{1}{h} \left\{ \mathbb{E} [N_v(X_{h,i})^2] - \mathbb{E} [N_v(X_{h,i})]^2 \right\}$$

and the estimator

$$\hat{\sigma}_v^2 = \frac{1}{h} \left\{ \frac{1}{n} \sum_{i=1}^n [N_v(X_{h,i})^2] - \left[ \frac{1}{n} \sum_{i=1}^n N_v(X_{h,i}) \right]^2 \right\}.$$

The statistic of interest for the generic Edgeworth expansion is, for  $1 \leq w \leq v \leq 2$ ,

$$T_{v,w} := \frac{\sqrt{nh}(\hat{f}_v - f)}{\hat{\sigma}_w}.$$

In this notation,

$$T_{\text{us}} = T_{1,1}, \quad T_{\text{bc}} = T_{2,1}, \quad \text{and} \quad T_{\text{rbc}} = T_{2,2}.$$

### S.I.1.2 Edgeworth Expansion Terms

The scaled bias is  $\eta_v = \sqrt{nh}(\mathbb{E}[\hat{f}_v] - f)$ . The Standard Normal distribution and density functions are  $\Phi(z)$  and  $\phi(z)$ , respectively.

The Edgeworth expansion for the distribution of  $T_{v,w}$  will consist of polynomials with coefficients that depend on moments of the kernel(s). To this end, continuing with the generic notation, for nonnegative integers  $j, k, p$ , define

$$\gamma_{v,p} = h^{-1} \mathbb{E} [N_v(X_{h,i})^p], \quad \Delta_{v,j} = \frac{1}{s} \sum_{i=1}^n \left\{ N_v(X_{h,i})^j - \mathbb{E} [N_v(X_{h,i})^j] \right\},$$

and

$$\nu_{v,w}(j, k, p) = \frac{1}{h} \mathbb{E} \left[ (N_v(X_{h,i}) - \mathbb{E} [N_v(X_{h,i})])^j (N_w(X_{h,i})^p - \mathbb{E} [N_w(X_{h,i})^p])^k \right].$$

We abbreviate  $\nu_{v,w}(j, 0, p) = \nu_v(j)$ .

To expand the distribution function, additional polynomials are needed beyond those used in the main text for coverage error. These are

$$\begin{aligned} p_{v,w}^{(1)}(z) &= \phi(z) \sigma_w^{-3} [\nu_{v,w}(1, 1, 2) z^2 / 2 - \nu_v(3)(z^2 - 1)/6], \\ p_{v,w}^{(2)}(z) &= -\phi(z) \sigma_w^{-3} \mathbb{E}[\hat{f}_w] \nu_{v,w}(1, 1, 1) z^2, \quad \text{and} \quad p_{v,w}^{(3)}(z) = \phi(z) \sigma_w^{-1}. \end{aligned}$$



Next, recall from the main text the polynomials used in *coverage error* expansions, here with an explicit argument for a generic quantile  $z$  rather than the specific  $z_{\alpha/2}$ :

$$\begin{aligned} q_1(z; K) &= \vartheta_{K,2}^{-2} \vartheta_{K,4}(z^3 - 3z)/6 - \vartheta_{K,2}^{-3} \vartheta_{K,3}^2 [2z^3/3 + (z^5 - 10z^3 + 15z)/9], \\ q_2(z; K) &= -\vartheta_{K,2}^{-1}(z), \quad \text{and} \quad q_3(z; K) = \vartheta_{K,2}^{-2} \vartheta_{K,3}(2z^3/3). \end{aligned}$$

The corresponding polynomials for expansions of the *distribution function* are

$$q_{v,w}^{(k)}(z) = \frac{1}{2} \frac{\phi(z)}{f} q_k(z; N_w), \quad k = 1, 2, 3.$$

Finally, the precise forms of  $\Omega_1$  and  $\Omega_2$  are:

$$\Omega_1 = -2 \frac{\mu_{K,\hat{k}}}{\nu_1(2)} \left\{ \int f(x - uh) K(u) L^{(\hat{k})}(u\rho) du - b \int f(x - uh) K(u) du \int f(x - ub) L^{(\hat{k})}(u) du \right\}$$

and  $\Omega_2 = \mu_{K,\hat{k}}^2 \vartheta_{K,2}^{-2} \vartheta_{L^{(\hat{k})},2}$ . These only appear for  $T_{bc}$ , and so are not indexed by  $\{v, w\}$ .

All these are discussed in Section S.I.6.

## S.I.2 Details of practical implementation

We maintain  $\ell = 2$  and recommend  $\hat{k} = 2$ . For the kernels  $K$  and  $L$ , we recommend either the second order minimum variance (to minimize interval length) or the MSE-optimal kernels; see Sections S.I.2.3 and S.I.4.2. In the next two subsections we discuss choice of  $h$  and  $\rho$ .

As argued below in Section S.I.2.3, we shall maintain  $\rho = 1$ . In the main text we give a direct plug-in (DPI) rule to implement the coverage-error optimal bandwidth. Here we give complete details for this procedure as well as document a second practical choice, based on a rule-of-thumb (ROT) strategy. Both choices yield the optimal coverage error decay rate of  $n^{-(\hat{k}+2)/(1+(\hat{k}+2))}$ .

All our methods are implemented in software available from the authors' websites and via the R package `nprobust` available at <https://cran.r-project.org/package=nprobust>.

**Remark 1** (Undercoverage of  $I_{us}(h_{mse}^*)$ ). It is possible not only to show that  $I_{us}(h_{mse}^*)$  asymptotically undercovers (see Hall and Horowitz (2013) for discussion in the regression context) but also to quantify precisely the coverage. To do so, write  $T_{us} = \sqrt{nh}(\hat{f} - \mathbb{E}[\hat{f}])/\hat{\sigma}_{us} + \eta_{us}/\hat{\sigma}_{us}$ , where the first term will be asymptotically standard Normal and the second will be a nonvanishing bias. To characterize the bias, recall from Eqn. (10) and Section S.I.1 that  $\eta_{us} = \sqrt{nh}h^{\hat{k}}[\mu_{K,\hat{k}}f^{(\hat{k})} + o(1)]$  and  $\hat{\sigma}^2 = \vartheta_{K,2}f[1 + o_P(1)]$ . Therefore, plugging in  $(h_{mse}^*)^{1+2\hat{k}} =$

$\vartheta_{K,2}f(\mu_{K,\hat{k}}f^{(\hat{k})})^{-2}/n$  shows that  $\eta_{\text{us}}/\hat{\sigma}_{\text{us}} = 1 + o_P(1)$ , whence  $T_{\text{us}}(h_{\text{mse}}^*) \rightarrow_d \mathcal{N}(1, 1)$ . For example, if  $\alpha = 0.05$ ,  $\mathbb{P}[f \in I_{\text{us}}(h_{\text{mse}}^*)] \approx 0.83$ . ■

### S.I.2.1 Bandwidth Choice: Rule-of-Thumb (ROT)

Motivated by the fact that estimating  $\hat{H}_{\text{dpi}}$  might be difficult in practice, while data-driven MSE-optimal bandwidth selectors are readily-available, the ROT bandwidth choice is to simply rescale any feasible MSE-optimal bandwidth  $\hat{h}_{\text{mse}}$  to yield optimal coverage error decay rates (but sub-optimal constants):

$$\hat{h}_{\text{rot}} = \hat{h}_{\text{mse}} n^{-(\hat{k}-2)/((1+2\hat{k})(\hat{k}+3))}.$$

When  $\hat{k} = 2$ ,  $\hat{h}_{\text{rot}} = \hat{h}_{\text{mse}}$ , which is optimal (in rates) as discussed previously.

**Remark 2** (Integrated Coverage Error). A closer analogue of the [Silverman \(1986\)](#) rule of thumb, which uses the integrated MSE, would be to integrate the coverage error over the point of evaluation  $x$ . For point estimation, this approach has some practical benefits. However, in the present setting note that  $\int f^{(\hat{k})}(x)dx = 0$ , removing the third term (of order  $h^{\hat{k}}$ ) entirely and thus, for any given point  $x$ , yields a lower quality approximation. ■

### S.I.2.2 Bandwidth Choice: Direct Plug-In (DPI)

To detail the direct plug-in (DPI) rule from the main text, it is useful to first simplify the problem. Recall from the main text that the optimal choice is  $h_{\text{rbc}}^* = H_{\text{rbc}}^*(\rho)n^{-1/(\hat{k}+3)}$ , where

$$\begin{aligned} H_{\text{rbc}}^*(K, L, \bar{\rho}) = \arg \min_H & \left| H^{-1} q_1(M_{\bar{\rho}}) + H^{1+2(\hat{k}+2)} (f^{(\hat{k}+2)})^2 (\mu_{K,\hat{k}+2} + \bar{\rho}^{-2} \mu_{K,\hat{k}} \mu_{L,2})^2 q_2(M_{\bar{\rho}}) \right. \\ & \left. + H^{\hat{k}+2} f^{(\hat{k}+2)} (\mu_{K,\hat{k}+2} + \bar{\rho}^{-2} \mu_{K,\hat{k}} \mu_{L,2}) q_3(M_{\bar{\rho}}) \right|. \end{aligned}$$

With  $\ell = 2$  and  $\rho = 1$ , and using the definitions of  $q_k(M_1)$ ,  $k = 1, 2, 3$ , from the main text or Section [S.I.1.2](#), this simplifies to:

$$\begin{aligned} H_{\text{rbc}}^*(K, L, 1) = \arg \min_H & \left| H^{-1} \left\{ \vartheta_{M,4} \frac{z^2 - 3}{6} - \vartheta_{M,3}^2 \frac{z^4 - 4z^2 + 15}{9} \right\} \right. \\ & - H^{1+2(\hat{k}+2)} \left\{ (f^{(\hat{k}+2)})^2 (\mu_{K,\hat{k}+2} + \mu_{K,\hat{k}} \mu_{L,2})^2 \vartheta_{M,2} \right\} \\ & \left. + H^{\hat{k}+2} \left\{ f^{(\hat{k}+2)} (\mu_{K,\hat{k}+2} + \mu_{K,\hat{k}} \mu_{L,2}) \vartheta_{M,3} \frac{2z^2}{3} \right\} \right|, \end{aligned}$$

where  $z = z_{\alpha/2}$  the appropriate upper quantile of the Normal distribution. However,  $H_{\text{rbc}}^*(\rho)$  still depends on the unknown density through  $f^{(\kappa+2)}$ .

Our recommendation is a DPI rule of order one, which uses a pilot bandwidth to estimate  $f^{(\kappa+2)}$  consistently. A simple and easy to implement choice is the MSE-optimal bandwidth appropriate to estimating  $f^{(\kappa+2)}$ , say  $h_{\kappa+2,\text{mse}}^*$ , which is different from  $h_{\text{mse}}^*$  for the level of the function; see e.g., [Wand and Jones \(1995\)](#). Let us denote a feasible MSE-optimal pilot bandwidth by  $\hat{h}_{\kappa+2,\text{mse}}$ . Then we have:

$$\begin{aligned} \hat{H}_{\text{dpi}}(K, L, 1) = \arg \min_H & \left| H^{-1} \left\{ \vartheta_{M,4} \frac{z^2 - 3}{6} - \vartheta_{M,3}^2 \frac{z^4 - 4z^2 + 15}{9} \right\} \right. \\ & - H^{1+2(\kappa+2)} \left\{ \hat{f}^{(\kappa+2)}(x; \hat{h}_{\kappa+2,\text{mse}})^2 (\mu_{K,\kappa+2} + \mu_{K,\kappa} \mu_{L,2})^2 \vartheta_{M,2} \right\} \\ & \left. + H^{\kappa+2} \left\{ \hat{f}^{(\kappa+2)}(x; \hat{h}_{\kappa+2,\text{mse}}) (\mu_{K,\kappa+2} + \mu_{K,\kappa} \mu_{L,2}) \vartheta_{M,3} \frac{2z^2}{3} \right\} \right|. \end{aligned}$$

This is now easily solved numerically (see note below). Further, if  $\kappa = 2$ , the most common case in practice, and  $K$  and  $L$  are either the respective second order minimum variance or MSE-optimal kernels (Sections [S.I.2.3](#) and [S.I.4.2](#)), then the above may be simplified to:

$$\begin{aligned} \hat{H}_{\text{dpi}}(M, 1) = \arg \min_H & \left| H^{-1} \left\{ \vartheta_{M,4} \frac{z^2 - 3}{6} - \vartheta_{M,3}^2 \frac{z^4 - 4z^2 + 15}{9} \right\} \right. \\ & - H^9 \left\{ \hat{f}^{(4)}(x; \hat{h}_{\kappa+2,\text{mse}})^2 \mu_{M,4}^2 \vartheta_{M,2} \right\} \\ & \left. + H^4 \left\{ \hat{f}^{(4)}(x; \hat{h}_{\kappa+2,\text{mse}}) \mu_{M,4} \vartheta_{M,3} \frac{2z^2}{3} \right\} \right|. \end{aligned}$$

Continuing with  $\kappa = 2$ , a second option is a DPI rule of order zero, which uses a reference model to build the rule of thumb, more akin to [Silverman \(1986\)](#). Using the Normal distribution, so that  $f(x) = \phi(x)$  and derivatives have known form, we obtain:

$$\begin{aligned} \hat{H}_{\text{dpi}}(M, 1) = \arg \min_H & \left| H^{-1} \left\{ \vartheta_{M,4} \frac{z^2 - 3}{6} - \vartheta_{M,3}^2 \frac{z^4 - 4z^2 + 15}{9} \right\} \right. \\ & - H^9 \left\{ [(\tilde{x}^4 - 6\tilde{x}^2 + 3) \phi(\tilde{x})]^2 \mu_{M,4}^2 \vartheta_{M,2} \right\} \\ & \left. + H^4 \left\{ (\tilde{x}^4 - 6\tilde{x}^2 + 3) \phi(\tilde{x}) \mu_{M,4} \vartheta_{M,3} \frac{2z^2}{3} \right\} \right| \end{aligned}$$

where  $\tilde{x} = (x - \hat{\mu})/\hat{\sigma}_X$  is the point of interest centered and scaled.

**Remark 3** (Notes on computation). When numerically solving the above minimization prob-

lems, computation will be greatly sped up by squaring the objective function. ■

### S.I.2.3 Choice of $\rho$

First, we expand on the argument that  $\rho$  should be bounded and positive. Intuitively, the standard errors  $\hat{\sigma}_{\text{rbc}}^2$  control variance up to order  $(nh)^{-1}$ , while letting  $b \rightarrow 0$  faster removes more bias. If  $b$  vanishes too fast, the variance is no longer controlled. Setting  $\bar{\rho} \in (0, \infty)$  balances these two. Let us simplify the discussion by taking  $\ell = 2$ , reflecting the widespread use of symmetric kernels. This does not affect the conclusions in any conceptual way, but considerably simplifies the notation. With this choice, Eqn. (9) yields the tidy expression

$$\eta_{\text{bc}} = \sqrt{nh} h^{\hat{k}+2} f^{(\hat{k}+2)} (\mu_{K,\hat{k}+2} - \rho^{-2} \mu_{K,\hat{k}} \mu_{L,2}) \{1 + o(1)\}.$$

Choice of  $\ell$  and  $b$  (or  $\rho$ ) cannot reduce the first term, which represents  $\mathbb{E}[\hat{f}] - f - B_f$ , and further, if  $\bar{\rho} = \infty$ , the bias rate is not improved, but the variance is inflated beyond order  $(nh)^{-1}$ . On the other hand, if  $\bar{\rho} = 0$ , then not only is a delicate choice of  $b$  needed, but  $\ell > 2$  is required, else the second term above dominates  $\eta_{\text{bc}}$ , and the full power of the variance correction is not exploited; that is, more bias may be removed without inflating the variance rate. Hall (1992b, p. 682) remarked that if  $\mathbb{E}[\hat{f}] - f - B_f$  is (part of) the leading bias term, then “explicit bias correction [...] is even less attractive relative to undersmoothing.” We show that, on the contrary, when using our proposed Studentization, it is optimal that  $\mathbb{E}[\hat{f}] - f - B_f$  is (part of) the dominant bias term. This reasoning is not an artifact of choosing  $\hat{k}$  even and  $\ell = 2$ , but in other cases  $\rho \rightarrow 0$  can be optimal if the convergence is sufficiently slow to equalize the two bias terms.

The following result which makes the above intuition precise.

**Corollary 6** (Robust bias correction:  $\rho \rightarrow 0$ ). *Let the conditions of Theorem 3(c) hold, with  $\bar{\rho} = 0$ , and fix  $\ell = 2$  and  $\hat{k} \leq S - 2$ . Then*

$$\begin{aligned} \mathbb{P}[f \in I_{\text{rbc}}] = 1 - \alpha + & \left\{ \frac{1}{nh} q_1(K) + nh^{1+2(\hat{k}+2)} (f^{(\hat{k}+2)})^2 (\mu_{K,\hat{k}+2}^2 + \rho^{-4} \mu_{K,\hat{k}}^2 \mu_{L,2}^2) q_2(K) \right. \\ & \left. + h^{\hat{k}+2} f^{(\hat{k}+2)} (\mu_{K,\hat{k}+2} + \rho^{-2} \mu_{K,\hat{k}} \mu_{L,2}) q_3(K) \right\} \frac{\phi(z_{\frac{\alpha}{2}})}{f} \{1 + o(1)\} \end{aligned}$$

By virtue of our new studentization, the leading variance remains order  $(nh)^{-1}$  and the problematic correlation terms are absent, however by forcing  $\rho \rightarrow 0$ , the  $\rho^{-2}$  terms of  $\eta_{\text{bc}}$  are dominant (the bias of  $\hat{B}_f$ ), and in light of our results, unnecessarily inflated. This verifies that  $\bar{\rho} = 0$  or  $\infty$  will be suboptimal.

We thus restrict to bounded and positive,  $\rho$ . Therefore,  $\rho$  impacts only the shape of the “kernel”  $M_\rho(u) = K(u) - \rho^{1+k} L^{(k)}(\rho u) \mu_{K,k}$ , and hence the choice of  $\rho$  depends on what properties the user desires for the kernel. It happens that  $\rho = 1$  has good theoretical properties and performs very well numerically (see Section S.I.8). As a result, from the practitioner’s point of view, choice of  $\rho$  (or  $b$ ) is completely automatic.

To see the optimality of  $\rho = 1$ , consider two cogent and well-studied possibilities: finding the kernel shape to minimize (i) interval length and (ii) MSE. The following optimal shapes are derived by Gasser et al. (1985) and references therein. Given the above results, we set  $k = 2$ . Indeed, the optimality properties here do not extend to higher order kernels.

Minimizing interval length is (asymptotically) equivalent to finding the minimum variance fourth-order kernel, as  $\sigma_{\text{rbc}}^2 \rightarrow f \vartheta_{M,2}$ . Perhaps surprisingly, choosing  $K$  and  $L^{(2)}$  to be the second-order minimum variance kernels for estimating  $f$  and  $f^{(2)}$  respectively, yields an  $M_1(u)$  that is exactly the minimum variance kernel. The fourth order minimum variance kernel for estimating  $f$  is  $K_{\text{mv}}(u) = (3/8)(-5u^2 + 3)$ , which is identical to  $M_1(u)$  when  $K$  is the uniform kernel and  $L^{(2)} = (15/4)(3u^2 - 1)$ , the minimum variance kernels for  $f$  and  $f^{(2)}$  respectively.

The result is similar for minimizing MSE: choosing  $K$  and  $L^{(2)}$  to be the MSE-optimal kernels for their respective point estimation problems yields an MSE-optimal  $M_1(u)$ . The optimal fourth order kernel is  $K_{\text{mse}}(u) = (15/32)(7u^4 - 10u^2 + 3)$ , and the respective second-order MSE optimal kernels are  $K(u) = (3/4)(1 - u^2)$  and  $L^{(2)}(u) = (105/16)(6u^2 - 5u^4 - 1)$ . A practitioner might use the MSE-optimal kernels (along with  $h_{\text{mse}}^*$ ) to obtain the best possible point estimate. Our results then give an accompanying measure of uncertainty that both has correct coverage and the attractive feature of using the same effective sample.

In Section S.I.4.2 we numerically compare several kernel shapes, focusing on: (i) interval length, measured by  $\vartheta_{M,2}$ , (ii) bias, given by  $\tilde{\mu}_{M,4}$ , and (iii) the associated MSE, given by  $(\vartheta_{M,2}^8 \tilde{\mu}_{M,4}^2)^{1/9}$ . These results, and the discussion above, give the foundations for our recommendation of  $\rho = 1$ , which delivers an easy-to-implement, fully automatic choice for implementing robust bias-correction that performs well numerically, as in Section S.I.8.

**Remark 4** (Coverage Error Optimal Kernels). Our results hint at a third notion of optimal kernel shape: minimizing coverage error. This kernel, for a fixed order  $k$ , would minimize the constants in Corollary 1 of the main text. In that result,  $h$  is chosen to optimize the rate and the constant  $H_{\text{us}}^*$  gives the minimum for a fixed kernel  $K$ . A step further would be to view  $H_{\text{us}}^*$  as a function of  $K$ , and optimizing. To our knowledge, such a derivation has not been done and may be of interest. ■

### S.I.3 Assumptions

Copied directly from the main text (see discussion there), the following assumptions are sufficient for our results.

**Assumption S.I.3.1** (Data-generating process).  *$\{X_1, \dots, X_n\}$  is a random sample with an absolutely continuous distribution with Lebesgue density  $f$ . In a neighborhood of  $x$ ,  $f > 0$ ,  $f$  is  $S$ -times continuously differentiable with bounded derivatives  $f^{(s)}$ ,  $s = 1, 2, \dots, S$ , and  $f^{(S)}$  is Hölder continuous with exponent  $\varsigma$ .*

**Assumption S.I.3.2** (Kernels). *The kernels  $K$  and  $L$  are bounded, even functions with support  $[-1, 1]$ , and are of order  $k \geq 2$  and  $\ell \geq 2$ , respectively, where  $k$  and  $\ell$  are even integers. That is,  $\mu_{K,0} = 1$ ,  $\mu_{K,k} = 0$  for  $1 \leq k < k$ , and  $\mu_{K,k} \neq 0$  and bounded, and similarly for  $\mu_{L,k}$  with  $\ell$  in place of  $k$ . Further,  $L$  is  $k$ -times continuously differentiable. For all integers  $k$  and  $l$  such that  $k + l = k - 1$ ,  $f^{(k)}(x_0)L^{(l)}((x_0 - x)/b) = 0$  for  $x_0$  in the boundary of the support.*

It will cause no confusion (as the notations never occur in the same place), but in the course of proofs we will frequently write  $s = \sqrt{nh}$ .

**Assumption S.I.3.3** (Cramér's Condition). *For each  $\xi > 0$  and all sufficiently small  $h$*

$$\sup_{t \in \mathbb{R}^2, t_1^2 + t_2^2 > \xi} \left| \int \exp\{i(t_1 M(u) + t_2 M(u)^2)\} f(x - uh) du \right| \leq 1 - C(x, \xi)h,$$

where  $C(x, \xi) > 0$  is a fixed constant and  $i = \sqrt{-1}$ .

**Remark 5** (Sufficient Conditions for Cramér's Condition). Assumption S.I.3.3 is a high level condition, but one that is fairly mild. Hall (1991) provides a primitive condition for Assumption S.I.3.3 and Lemma 4.1 in that paper verifies that Assumption S.I.3.3 is implied. Hall (1992a) and Hall (1992b) assume the same primitive condition. This condition is as follows. On their compact support, assumed here to be  $[-1, 1]$ , there exists a partition  $-1 = a_0 < a_1 < \dots < a_m = 1$ , such that on each  $(a_{j-1}, a_j)$ ,  $K$  and  $M$  are differentiable, with bounded, strictly monotone derivatives.

This condition is met for many kernels, with perhaps the only exception of practical importance being the uniform kernel. As Hall (1991) describes, it is possible to prove the Edgeworth expansion for the uniform kernel using different methods than we use in below. The uniform kernel is also ruled out for local polynomial regression, see Remark 9. ■

## S.I.4 Bias

This section accomplishes three things. First, we first carefully derive the bias of the initial estimator and the bias correction. Second, we explicate the properties of the induced kernel  $M_\rho$  in terms of bias reduction and how exactly this kernel is “higher-order”. Finally, we examine two other methods of bias reduction: (i) estimating the derivatives without using derivatives of kernels (Singh, 1977), and (ii) the generalized jackknife approach (Schucany and Sommers, 1977). Further methods are discussed and compared by Jones and Signorini (1997). The message from both alternative methods echoes our main message: it is important to account for any bias correction when doing inference, i.e., to avoid the mismatch present in  $T_{bc}$ .

### S.I.4.1 Precise Bias Calculations

Recall that the biases of the two estimators are as follows:

$$\mathbb{E}[\hat{f}] - f = \begin{cases} h^{\hat{k}} f^{(\hat{k})} \mu_{K,\hat{k}} + h^{\hat{k}+2} f^{(\hat{k}+2)} \mu_{K,\hat{k}+2} + o(h^{\hat{k}+2}) & \text{if } \hat{k} \leq S-2 \\ h^{\hat{k}} f^{(\hat{k})} \mu_{K,\hat{k}} + O(h^{S+\varsigma}) & \text{if } \hat{k} \in \{S-1, S\} \\ 0 + O(h^{S+\varsigma}) & \text{if } \hat{k} > S \end{cases} \quad (10)$$

and

$$\mathbb{E}[\hat{f} - \hat{B}_f] - f = \begin{cases} h^{\hat{k}+2} f^{(\hat{k}+2)} \mu_{K,\hat{k}+2} + h^{\hat{k}} b^{\ell} f^{(\hat{k}+\ell)} \mu_{K,\hat{k}} \mu_{L,\ell} + o(h^{\hat{k}+2} + h^{\hat{k}} b^{\ell}) & \text{if } \hat{k} + s \leq S \\ h^{\hat{k}+2} f^{(\hat{k}+2)} \mu_{K,\hat{k}+2} + O(h^{\hat{k}} b^{S-\hat{k}+\varsigma}) + o(h^{\hat{k}+2}) & \text{if } 2 \leq S - \hat{k} < s \\ O(h^{S+\varsigma}) + O(h^{\hat{k}} b^{S-\hat{k}+\varsigma}) & \text{if } \hat{k} \in \{S-1, S\} \\ O(h^{S+\varsigma}) + O(h^{\hat{k}} b^{S-\hat{k}}) & \text{if } \hat{k} > S. \end{cases} \quad (11)$$

The following Lemma gives a rigorous proof of these statements.

**Lemma 1.** *Under Assumptions S.I.3.1 and S.I.3.2, Equations (10) and (11) hold.*

*Proof.* To show Eqn. (10), begin with the change of variables and the Taylor expansion

$$\begin{aligned} \mathbb{E}[\hat{f}] &= h^{-1} \int K(X_{h,i}) f(X_i) dX_i = \int K(u) f(x - uh) du \\ &= \sum_{k=0}^S \left\{ (-h)^k f^{(k)}(x) \int u^k K(u) du / k! \right\} + (-h)^S \int u^S K(u) (f^{(S)}(\bar{x}) - f^{(S)}(x)) du. \end{aligned}$$

where  $\bar{x} \in [x, x - uh]$ . By the Hölder condition of Assumption S.I.3.1, the final term is  $O(h^{S+\varsigma})$ . If  $k > S$ , then all  $\int u^k K(u) du = 0$ , and only this remainder is left. In all other cases,  $h^k f^{(k)}(x) \mu_{K,k}$  is the first nonzero term of the summation, and hence the leading bias term. Further, by virtue of  $k$  being even and  $K$  symmetric,  $\int u^{k+1} K(u) du = 0$ , leaving only  $O(h^{S+\varsigma})$  when  $k = S-1$ , and otherwise, when  $k \leq S-2$ , leaving  $h^{k+2} f^{(k+2)}(x) \mu_{K,k+2} + o(h^{k+2})$ . This completes the proof of Eqn. (10).

To establish Eqn. (11), first write

$$\mathbb{E}[\hat{f} - \hat{B}_f] - f = \mathbb{E}[\hat{f} - f - B_f] + \mathbb{E}[B_f - \hat{B}_f],$$

where  $B_f$  follows the convention of being identically zero if  $k > S$ . The first portion is characterized by rearranging Eqn. (10), so it remains to examine the second term. Let  $\tilde{k} = k \vee S$ . By repeated integration by parts, using the boundary conditions of Assumption S.I.3.2:

$$\begin{aligned} \mathbb{E}[\hat{f}^{(\tilde{k})}] &= \frac{1}{b^{1+\tilde{k}}} \int L^{(\tilde{k})}(X_{b,i}) f(X_i) dX_i \\ &= -\frac{1}{b^{1+(\tilde{k}-1)}} L^{(\tilde{k}-1)}(X_{b,i}) f(X_i) \Big|_x + \frac{1}{b^{1+(\tilde{k}-1)}} \int L^{(\tilde{k}-1)}(X_{b,i}) f^{(1)}(X_i) dX_i \\ &= 0 + \frac{1}{b^{1+(\tilde{k}-1)}} \int L^{(\tilde{k}-1)}(X_{b,i}) f^{(1)}(X_i) dX_i \\ &= -\frac{1}{b^{1+(\tilde{k}-2)}} L^{(\tilde{k}-2)}(X_{b,i}) f^{(1)}(X_i) + \frac{1}{b^{1+(\tilde{k}-2)}} \int L^{(\tilde{k}-2)}(X_{b,i}) f^{(2)}(X_i) dX_i \\ &\quad \vdots \\ &= \frac{1}{b^{1+(\tilde{k}-\tilde{k})}} \int L^{(\tilde{k}-\tilde{k})}(X_{b,i}) f^{(\tilde{k})}(X_i) dX_i \\ &= \frac{1}{b^{\tilde{k}-\tilde{k}}} \int L^{(\tilde{k}-\tilde{k})}(u) f^{(\tilde{k})}(x - ub) du, \end{aligned}$$

where the last line follows by a change of variables. We now proceed separately for each case delineated in (11), from top to bottom. For  $k > S$ , no reduction is possible, and the final line above is  $O(b^{S-k})$ , and with  $B_f = 0$ , we have  $\mathbb{E}[B_f - \hat{B}_f] = 0 - h^k \mu_{K,k} \mathbb{E}[\hat{f}^{(k)}] = O(h^k b^{S-k})$ , as shown. For  $k \leq S$ , by a Taylor expansion, the final line displayed above becomes

$$\sum_{k=\tilde{k}}^S \{b^{k-\tilde{k}} f^{(k)}(x) \mu_{L,k-\tilde{k}}\} + b^{S-\tilde{k}} \int u^{S-\tilde{k}} L(u) (f^{(S)}(\bar{x}) - f^{(S)}(x)) du.$$

The second term above is  $O(b^{S-k+\varsigma})$  in all cases, and  $\mu_{L,0} = 1$ , which yields  $\mathbb{E}[\hat{f}^{(\tilde{k})}] = f^{(\tilde{k})} + O(b^{S-\tilde{k}+\varsigma})$  for  $k \in \{S-1, S\}$ , using  $\mu_{L,1} = 0$  in the former case. Next, if  $k + \ell \leq S$ , the above becomes  $\mathbb{E}[\hat{f}^{(\tilde{k})}] = f^{(\tilde{k})} + b^\ell f^{(k+\ell)} \mu_{L,\ell} + o(b^\ell)$ , as  $\mu_{L,k} = 0$  for  $1 < k < \ell$ , whereas if  $k + \ell > S$ ,



the remainder terms can not be characterized, leaving  $\mathbb{E}[\hat{f}^{(\hat{k})}] = f^{(\hat{k})} + O(b^{S-\hat{k}+\varsigma})$ . Plugging any of these results into  $\mathbb{E}[B_f - \hat{B}_f] = h^{\hat{k}} \mu_{K,\hat{k}}(f^{(\hat{k})} - \mathbb{E}[\hat{f}^{(\hat{k})}])$  completes the demonstration of Eqn. (11).  $\square$

### S.I.4.2 Properties of the kernel $M_\rho(\cdot)$

As made precise below,  $M_\rho$  is a higher-order kernel. The choices of  $K$ ,  $L$ , and  $\rho$  determine the shape of  $M_\rho$ , which in turn effects the variance and bias constants. In standard kernel analyses, these constants are used to determine optimal kernel shapes for certain problems (see Gasser et al. (1985) and references therein). For several choices of  $K$ ,  $L$ , and  $\rho$ , Table S.I.1 shows numerical results for the various constants of the induced kernel  $M_\rho$ . The table includes (i) the variance, given by  $\vartheta_{M,2}$  and relevant for interval length, (ii) a measure of bias given by  $\tilde{\mu}_{M,4}$ , and finally (iii) the resulting mean square error constant,  $[\vartheta_{M,2}^8 \tilde{\mu}_{M,4}^2]^{1/9}$  ( $\tilde{\mu}_{M,4} = (k!)(-1)^k \mu_{M,4}$ ). These specific constants are due to  $M_\rho$  being a fourth order kernel, as discussed next, and would otherwise remain conceptually the same but rely on different moments. A more general, but more cumbersome procedure would be to choose  $\rho$  numerically to minimize some notation of distance (e.g.,  $L_2$ ) between the resulting kernel  $M_\rho$  and the optimal kernel shape already available in the literature. However, using  $\rho = 1$  as a simple rule-of-thumb exhibits very little lost performance, as shown in the Table and discussed in the paper.

It is worthwhile to make precise the sense in which the  $n$ -varying “kernel”  $M_\rho(\cdot)$  of Eqn. (9) is a higher-order kernel. Comparing Equations (10) and (11) shows exactly what is meant by this statement: the bias rate attained agrees with a standard estimate using a kernel of order  $\hat{k} + 2$  (if  $\bar{\rho} > 0$ ), as  $\ell \geq 2$ . For example, if  $\hat{k} = \ell = 2$  and  $\bar{\rho} > 0$ , then  $M_{\bar{\rho}}(\cdot)$  behaves as a fourth-order kernel in terms of bias reduction.

However, it is not true in general that  $M(\cdot)$  is a higher-order kernel in the sense that its moments below  $\hat{k} + 2$  are zero. That is, for any  $k < \hat{k}$ , by the change of variables  $w = \rho u$ ,

$$\begin{aligned} \int_{-1}^1 u^k M(u) du &= \int_{-1}^1 u^k K(u) du - \rho^{1+\hat{k}} \mu_{K,\hat{k}} \int_{-1}^1 u^k L^{(\hat{k})}(\rho u) du \\ &= 0 - \rho^{1+\hat{k}} \mu_{K,\hat{k}} \rho^{-1-k} \int_{-\rho}^{\rho} w^k L^{(\hat{k})}(w) dw \\ &= 0 - \rho^{\hat{k}-k} \mu_{K,\hat{k}} \int_{-\rho}^{\rho} w^k L^{(\hat{k})}(w) dw. \end{aligned}$$

Now,  $L(u) = L(-u)$  implies that  $L^{(k)}(u) = (-1)^k L^{(k)}(-u)$ . Since  $\hat{k}$  is even,  $L^{(\hat{k})}(w)$  is symmetric, therefore if  $k$  is odd  $0 = \int_{-\rho}^{\rho} w^k L^{(\hat{k})}(w) dw$  for any  $\rho$ . But this fails for  $k$  even, even for  $\rho = 1$ , and hence  $\int_{-1}^1 u^k M(u) du \neq 0$ . For example, in the leading case of  $\hat{k} = \ell = 2$ ,

Table S.I.1: Numerical results for bias and variance constants of the induced higher-order kernel  $M$  for several choices of  $K$ ,  $L$ , and  $\rho$

Kernel $K$	Kernel $L^{(2)}$	$\rho = 0.5$				$\rho = 1$				$\rho = 1.5$			
		$\tilde{\mu}_{M,4}$	$\vartheta_{M,2}$	MSE		$\tilde{\mu}_{M,4}$	$\vartheta_{M,2}$	MSE		$\tilde{\mu}_{M,4}$	$\vartheta_{M,2}$	MSE	
Epanechnikov	$(105/16)(6u^2 - 5u^4 - 1)$	0.0690	0.6430	0.3729		-0.0476	1.2500	0.6199		-0.3643	5.5992	3.6944	
Uniform	$(105/16)(6u^2 - 5u^4 - 1)$	0.1722	0.5152	0.3752		-0.0222	1.4722	0.6052		-0.5500	11.5742	7.7202	
Biweight	$(105/16)(6u^2 - 5u^4 - 1)$	0.0357	0.7617	0.3744		-0.0476	1.2500	0.6199		-0.2738	3.9537	2.5448	
Triweight	$(105/16)(6u^2 - 5u^4 - 1)$	0.0210	0.8617	0.3715		-0.0438	1.2774	0.6202		-0.2197	3.2395	2.0300	
Tricube	$(105/16)(6u^2 - 5u^4 - 1)$	0.0335	0.7542	0.3658		-0.0506	1.2332	0.6207		-0.2786	3.9344	2.5436	
Cosine	$(105/16)(6u^2 - 5u^4 - 1)$	0.0629	0.6617	0.3747		-0.0476	1.2503	0.6199		-0.3475	5.2717	3.4651	
Epanechnikov	$(15/4)(3u^2 - 1)$	0.0643	0.6410	0.3660		-0.0857	1.1250	0.6432		-0.4929	4.1754	3.0440	
Uniform	$(15/4)(3u^2 - 1)$	0.1643	0.5098	0.3678		-0.0857	1.1250	0.6432		-0.7643	7.6191	5.7276	
Biweight	$(15/4)(3u^2 - 1)$	0.0323	0.7543	0.3630		-0.0748	1.1352	0.6291		-0.3656	3.0550	2.1579	
Triweight	$(15/4)(3u^2 - 1)$	0.0184	0.8517	0.3568		-0.0649	1.1631	0.6229		-0.2911	2.5444	1.7435	
Tricube	$(15/4)(3u^2 - 1)$	0.0300	0.7487	0.3547		-0.0780	1.1319	0.6333		-0.3712	3.0729	2.1764	
Cosine	$(15/4)(3u^2 - 1)$	0.0584	0.6583	0.3669		-0.0836	1.1254	0.6399		-0.4693	3.9510	2.8668	
Biweight	Biweight <sup>(2)</sup>	0.0323	0.7543	0.3630		-0.0748	1.1352	0.6291		-0.3656	3.0550	2.1579	
Tricube	Tricube <sup>(2)</sup>	0.0299	0.7516	0.3556		-0.0790	1.1993	0.6687		-0.3746	3.7063	2.5762	
Gaussian	Gaussian <sup>(2)</sup>	2.2500	0.3006	0.4113		-3.0000	0.4760	0.6599		-17.2500	1.3606	2.4758	

<sup>1</sup> As discussed in Section S.I.4.2,  $M_\rho$  behaves as a fourth order kernel in terms of bias reduction, but does not strictly fit within the class of kernels used in derivation of optimal kernel shapes. This explains the super-optimal behavior exhibited by some choices of  $K$ ,  $L$ , and  $\rho$ .

<sup>2</sup> The constants  $\tilde{\mu}_{M,4}$  and  $\vartheta_{M,2}$  measure bias and variance, respectively (the latter also being relevant for interval length). The MSE is measured by  $[\vartheta_{M,2}^8 \tilde{\mu}_{M,4}^2]^{1/9}$ , owing to  $M_\rho$  being a fourth-order kernel.

$\int_{-1}^1 u^2 M(u) du \neq 0$  in general, and so  $M(\cdot)$  is not a fourth-order kernel in the traditional sense.

Instead, the bias reduction is achieved differently. The proof of Lemma 1 makes explicit use of the structure imposed by estimating  $f^{(\kappa)}$  using the *derivative* of the kernel  $L(\cdot)$ . From a technical standpoint, an integration by parts argument shows how the properties of the kernel  $L(\cdot)$  (not the function  $L^{(\kappa)}(\cdot)$ ) are used to reduce bias. This argument *precedes* the Taylor expansion of  $f$ , and thus moments of  $M$  are never encountered and there is no requirement that they be zero. This approach is simple, intuitive, and leads to natural restrictions on the kernel  $L$ , and for this reason it is commonly employed in the literature and in practice (Hall, 1992b).

### S.I.4.3 Other Bias Reduction Methods

We now examine two other methods of bias reduction: (i) estimating the derivatives without using derivatives of kernels (Singh, 1977), and (ii) the generalized jackknife approach (Schucany and Sommers, 1977). Further methods are discussed and compared by Jones and Signorini (1997). Both methods are shown to be tightly connected to our results. Further, a more general message is that it is important to account for any bias correction when doing inference, i.e., to avoid the mismatch present in  $T_{bc}$ .

The first method, which dates at least to Singh (1977), is to introduce a class of kernel functions directly for derivative estimation, more closely following the standard notion of a higher-order kernel rather than using the derivative of a kernel to estimate the density derivative and proving bias reduction via integration by parts. Jones (1994) expands on this method and gives further references. This class of kernels is used in the derivation of optimal kernel shapes (for derivative estimation) by Gasser et al. (1985). It is worthwhile to show how this class of kernel achieves bias correction and how this approach fits into our Edgeworth expansions.

Consider estimating  $f^{(\kappa)}$  with

$$\tilde{f}^{(\kappa)}(x) = \frac{1}{nb^{1+\kappa}} \sum_{i=1}^n J(X_{b,i}),$$

for some kernel function  $J(\cdot)$ . Note well that  $J$  is generic, it need not itself be a derivative, but this is the only difference here. A direct Taylor expansion (i.e. without first integrating by parts) then gives

$$\mathbb{E}[\tilde{f}^{(\kappa)}] = b^{-\kappa} \sum_{k=0}^S b^k \mu_{J,k} f^{(k)} + O(b^{S+\varsigma}).$$

Thus, if  $J$  satisfies  $\mu_{J,k} = 0$  for  $k = 0, 1, \dots, k-1, k+1, k+2, \dots, k+(\ell-1)$ ,  $\mu_{J,k} = 1$ , and  $\mu_{J,k+\ell} \neq 0$ , and  $S$  is large enough then

$$\mathbb{E}[\tilde{f}^{(k)}] = f^{(k)} + b^\ell f^{(k+\ell)} \mu_{J,k+\ell} + o(b^\ell),$$

just as achieved by  $\hat{f}^{(k)}$  and exactly matching Eqn. (10). Note that  $\mu_{J,0} = 0$ , that is, the kernel  $J$  does not integrate to one. In the language of Gasser et al. (1985),  $J$  is a kernel of order  $(k, k+\ell)$ .

Given this result, bias correction can of course be performed using  $\tilde{f}^{(k)}(x)$  (based on  $J$ ) rather than  $\hat{f}^{(k)}$  (based on  $L^{(k)}$ ). Much will be the same: the structure of Eqn. (9) will hold with  $J$  in place of  $L^{(k)}$  and the results in Eqn. (11) are achieved with modifications to the constants (e.g., in the first line,  $\mu_{J,k+\ell}$  appears in place of  $\mu_{L,\ell}$ ). In either case, the same bias rates are attained. Our Edgeworth expansions will hold for this class under the obvious modifications to the notation and assumptions, and all the same conclusions are obtained.

When studying optimal kernel shapes, Gasser et al. (1985) actually further restrict the class, by placing a limit on the number of sign changes over the support of the kernel, which ensures that the MSE and variance minimization problems have well-defined solutions. Collectively, these differences in the kernel classes explain why it is possible to demonstrate “super-optimal” MSE and variance performance for certain choices of  $K$ ,  $L^{(k)}$ , and  $\rho$ , as in Table S.I.1.

A second alternative is the generalized jackknife method of Schucany and Sommers (1977), and expanded upon by Jones and Foster (1993). To simplify the notation and ease exposition, we describe this approach for second order kernels ( $k = 2$ ), but the method, and all the conclusions below, generalize fully. We thank an anonymous reviewer for encouraging us to include these details.

Begin with two estimators  $\hat{f}_1$  and  $\hat{f}_2$ , with (possibly different) bandwidths and second-order kernels  $h_j$  and  $K_j$ ,  $j = 1, 2$ ; thus Eqn. (10) gives

$$\mathbb{E}[\hat{f}_j] - f(x) = h_j^2 f^{(2)} \mu_{K_j,2} + o(h_j^2), \quad j = 1, 2.$$

Schucany and Sommers (1977) propose to estimate  $f$  with  $\hat{f}_{\text{GJ},R} := (\hat{f}_1 - R\hat{f}_2)/(1-R)$ , the bias of which is

$$\mathbb{E}[\hat{f}_{\text{GJ},R} - f] = \frac{f^{(2)}}{1-R} (h_1^2 \mu_{K_1,2} - R h_2^2 \mu_{K_2,2}) + o(h_1^2 + h_2^2).$$

Hence, setting  $R = (h_1^2 \mu_{K_1,2}) / (h_2^2 \mu_{K_2,2})$  renders the leading bias exactly zero. Moreover, if  $S \geq 4$ ,  $\hat{f}_{\text{GJ},R}$  has bias  $O(h_1^4 + h_2^4)$ ; behaving as a single estimator with  $k = 4$ . To put this in

context of our results, observe that with this choice of  $R$ , if we let  $\tilde{\rho} = h_1/h_2$ , then

$$\hat{f}_{\text{GJ},R} = \frac{1}{nh_1} \sum_{i=1}^n \tilde{M}\left(\frac{X_i - x}{h_1}\right), \quad M(u) = K_1(u) - \tilde{\rho}^{1+2} \left\{ \frac{K_2(\tilde{\rho}u) - \tilde{\rho}^{-1}K_1(u)}{\mu_{K_2,2}(1-R)} \right\} \mu_{K_1,2},$$

exactly matching Eqn. (9). Or equivalently,  $\hat{f}_{\text{GJ},R} = \hat{f}_1 - h_1^2 \tilde{f}^{(2)} \mu_{K_1,2}$ , for the derivative estimator

$$\tilde{f}^{(2)} = \frac{1}{nh_2^{1+2}} \sum_{i=1}^n \tilde{L}\left(\frac{X_i - x}{h_2}\right), \quad \tilde{L}(u) = \frac{K_2(u) - \tilde{\rho}^{-1}K_1(\tilde{\rho}^{-1}u)}{\mu_{K_2,2}(1-R)}.$$

Therefore, we can view  $\hat{f}_{\text{GJ},R}$  as a change in the kernel  $M(\cdot)$  or an explicit bias estimation described directly above with a specific choice of  $J(\cdot)$  (depending on  $\tilde{\rho}$  in either case). Again, Eqn. (9) holds exactly. Thus, our results cover the generalized jackknife method as well, and the same lessons apply.

Finally, we note that these bias correction methods can be applied to nonparametric regression as well, and local polynomial regression in particular, and that the same conclusions are found. We will not repeat this discussion however.

## S.I.5 First Order Properties

Here we briefly state the first-order properties of  $T_{\text{us}}$ ,  $T_{\text{bc}}$ , and  $T_{\text{rbc}}$ , using the common notation  $T_{v,w}$  defined in Section S.I.1. Recall that  $\eta_v = \sqrt{nh}(\mathbb{E}[\hat{f}_v] - f)$  is the scaled bias in either case. With this notation, we have the following result.

**Lemma 2.** *Let Assumptions S.I.3.1 and S.I.3.2 hold. Then if  $nh \rightarrow \infty$ ,  $\eta_v \rightarrow 0$ , and if  $v = 2$ ,  $\rho \rightarrow 0 + \bar{\rho} \mathbb{1}\{v = w\} < \infty$ , it holds that  $T_{v,w} \rightarrow_d \mathcal{N}(0, 1)$ .*

The conditions on  $h$  and  $b$  behind the generic assumption that the scaled bias vanishes can be read off of (10) and (11):  $T_{\text{us}}$  requires  $\sqrt{nh}h^k \rightarrow 0$  whereas  $T_{\text{bc}}$  and  $T_{\text{rbc}}$  require only  $\sqrt{nh}h^k(h^2 \vee b^\ell) \rightarrow 0$ , and thus accommodate  $\sqrt{nh}h^k \not\rightarrow 0$  or  $b \not\rightarrow 0$  (but not both). However, bias correction requires a choice of  $\rho = h/b$ . One easily finds that  $\mathbb{V}[\sqrt{nh}\hat{B}_f] = O(\rho^{1+2k})$ , whence  $\rho \rightarrow 0$  is required for  $T_{\text{bc}}$ . But  $T_{\text{rbc}}$  does not suffer from this requirement because of our proposed, new Studentization. From a first-order point of view, traditional bias correction allows for a larger class of sequences  $h$ , but requires a delicate choice of  $\rho$  (or  $b$ ), and Hall (1992b) shows that this constraint prevents  $T_{\text{bc}}$  from improving inference. Our novel standard errors remove these constraints, allowing for improvements in bias to carry over to improvements in inference. The fact that a wider range of bandwidths is allowed hints at

the robustness to tuning parameter choice discussed above and formalized by our Edgeworth expansions.

**Remark 6** ( $\rho \rightarrow \infty$ ).  $T_{\text{rbc}} \rightarrow_d \mathcal{N}(0, 1)$  will hold even for  $\bar{\rho} = \infty$ , under the even weaker bias rate restriction that  $\eta_{\text{bc}} = o(\rho^{1/2+\kappa})$ , provided  $nb \rightarrow \infty$ . In this case  $\hat{B}_f$  dominates the first-order approximation, but  $\sigma_{\text{rbc}}^2$  still accounts for the total variability. However there is no gain for inference: the bias properties can not be improved due to the second bias term  $(\mathbb{E}[\hat{f}] - f - B_f)$ , while variance can only be inflated. Thus, we restrict to bounded  $\bar{\rho}$ . Section S.I.2.3 has more discussion on the choice of  $\rho$ . ■

## S.I.6 Main Result: Edgeworth Expansion

Recall the generic notation:

$$T_{v,w} := \frac{\sqrt{nh}(\hat{f}_v - f)}{\hat{\sigma}_w},$$

for  $1 \leq w \leq v \leq 2$ . The Edgeworth expansion for the distribution of  $T_{v,w}$  will consist of polynomials with coefficients that depend on moments of the kernel(s). Additional polynomials are needed beyond those used in the main text for coverage error. These are:

$$\begin{aligned} p_{v,w}^{(1)}(z) &= \phi(z)\sigma_w^{-3}[\nu_{v,w}(1, 1, 2)z^2/2 - \nu_v(3)(z^2 - 1)/6], \\ p_{v,w}^{(2)}(z) &= -\phi(z)\sigma_w^{-3}\mathbb{E}[\hat{f}_w]\nu_{v,w}(1, 1, 1)z^2, \quad \text{and} \quad p_{v,w}^{(3)}(z) = \phi(z)\sigma_w^{-1}. \end{aligned}$$

The polynomials  $p_{v,w}^{(k)}$  are even, and hence cancel out of coverage probability expansions, but are used in the expansion of the distribution function itself (or equivalently, the coverage of a one-sided confidence interval).

Next, recall from the main text the polynomials used in *coverage error* expansions:

$$\begin{aligned} q_1(z; K) &= \vartheta_{K,2}^{-2}\vartheta_{K,4}(z^3 - 3z)/6 - \vartheta_{K,2}^{-3}\vartheta_{K,3}^2[2z^3/3 + (z^5 - 10z^3 + 15z)/9], \\ q_2(z; K) &= -\vartheta_{K,2}^{-1}(z), \quad \text{and} \quad q_3(z; K) = \vartheta_{K,2}^{-2}\vartheta_{K,3}(2z^3/3). \end{aligned}$$

The corresponding polynomials for expansions of the *distribution function* are

$$q_{v,w}^{(k)}(z) = \frac{1}{2} \frac{\phi(z)}{f} q_k(z; N_w), \quad k = 1, 2, 3.$$

As before, the  $q_{v,w}^{(k)}$  are odd and hence do not cancel when computing coverage: the  $q_k(z; N_w)$  in the main text are doubled for just this reason.

Note that, despite the notation,  $q_{v,w}^{(k)}(z)$  depends only on the “denominator” kernel  $N_w$ . The notation comes from the fact that when first computed, the terms which enter into the  $q_{v,w}^{(k)}(z)$  depend on both kernels, but the simplifications in Eqn. (16) reduce the dependence to  $N_w$ . This is because for undersmoothing and robust bias correction,  $v = w$ , and for traditional bias correction  $N_2 = M = K + o(1) = N_1 + o(1)$ , as  $\rho \rightarrow 0$  is assumed. Thus, when computing  $\vartheta_{M,q}$  the terms with the lowest powers of  $\rho$  will be retained. These can be found by expanding

$$\vartheta_{M,q} = \int (K(u) - \rho^{1+\kappa} \mu_{K,\kappa} L^{(\kappa)}(u))^q du = \sum_{j=0}^q \binom{q}{j} (-\mu_{K,\kappa} \rho^{1+\kappa})^{q-j} \int K(u)^j L^{(\kappa)}(\rho u)^{q-j} du,$$

and hence we can write  $\vartheta_{M,q} = \vartheta_{K,q} - \rho^{1+\kappa} q \mu_{K,\kappa} L^{(\kappa)}(0) \vartheta_{K,q-1} + O(h + \rho^{2+\kappa})$ . We can thus write  $q_j(z; M) = q_j(z; K) + o(1)$  in this case. If the expansions were carried out beyond terms of order  $(nh)^{-1} + (nh)^{-1/2} \eta_v + \eta_v^2 + \mathbb{1}\{v \neq w\} \rho^{1+2\kappa}$  this would not be the case.

Finally, for traditional bias correction, there are additional terms in the expansion (see discussion in the main text) representing the covariance of  $\hat{f}$  and  $\hat{B}_f$  (denoted by  $\Omega_1$ ) and the variance of  $\hat{B}_f$  ( $\Omega_2$ ). We now state their precise forms. These arise from the mismatch between the variance of the numerator of  $T_{bc}$  and the standardization used,  $\sigma_{us}^2$ , that is  $\sigma_{rbc}^2 / \sigma_{us}^2$  is given by

$$\frac{nh\mathbb{V}[\hat{f} - \hat{B}_f]}{nh\mathbb{V}[\hat{f}]} = \frac{nh\mathbb{V}[\hat{f}] - 2nh\mathbb{C}[\hat{f}, \hat{B}_f] + nh\mathbb{V}[\hat{B}_f]}{nh\mathbb{V}[\hat{f}]} = 1 - 2\frac{nh\mathbb{C}[\hat{f}, \hat{B}_f]}{nh\mathbb{V}[\hat{f}]} + \frac{nh\mathbb{V}[\hat{B}_f]}{nh\mathbb{V}[\hat{f}]}.$$

This makes clear that  $\Omega_1$  and  $\Omega_2$  are the constant portions of the last two terms. We have

$$-2\frac{nh\mathbb{C}[\hat{f}, \hat{B}_f]}{nh\mathbb{V}[\hat{f}]} = \rho^{1+\kappa} \Omega_1,$$

where

$$\Omega_1 = -2\frac{\mu_{K,\kappa}}{\nu_1(2)} \left\{ \int f(x - uh) K(u) L^{(\kappa)}(u\rho) du - b \int f(x - uh) K(u) du \int f(x - ub) L^{(\kappa)}(u) du \right\}.$$

Note  $\nu_1(2) = \sigma_{us}^2$ . Turning to  $\Omega_2$ , using the calculations in Section S.I.4.1 (recall  $\tilde{\kappa} = \kappa \vee S$ ), we find that

$$\frac{nh\mathbb{V}[\hat{B}_f]}{nh\mathbb{V}[\hat{f}]} = \rho^{1+2\kappa} \Omega_2 \quad \text{where} \quad \Omega_2 = \frac{\mu_{K,\kappa}^2}{\nu_1(2)} \left\{ \int f(x - ub) L^{(\kappa)}(u)^2 du - b^{1+2\tilde{\kappa}} \left( \int L^{(\kappa-\tilde{\kappa})}(u) f^{(\tilde{\kappa})}(x - ub) du \right) \right\}.$$

Fully simplifying would yield

$$\Omega_2 = \mu_{K,\hat{k}}^2 \vartheta_{K,2}^{-2} \vartheta_{L^{(\hat{k})},2},$$

which can be used in Theorem 3.

As a last piece of notation, define the scaled bias as  $\eta_v = \sqrt{nh}(\mathbb{E}[\hat{f}_v] - f)$ .

We can now state our generic Edgeworth expansion, from whence the coverage probability expansion results follow immediately.

**Theorem 3.** *Suppose Assumptions S.I.3.1, S.I.3.2, and S.I.3.3 hold,  $nh/\log(n) \rightarrow \infty$ ,  $\eta_v \rightarrow 0$ , and if  $v = 2$ ,  $\rho \rightarrow 0 + \bar{\rho}\mathbb{1}\{v = w\}$ . Then for*

$$\begin{aligned} F_{v,w}(z) = & \Phi(z) + \frac{1}{\sqrt{nh}}p_{v,w}^{(1)}(z) + \sqrt{\frac{h}{n}}p_{v,w}^{(2)}(z) + \eta_v p_{v,w}^{(3)}(z) + \frac{1}{nh}q_{v,w}^{(1)}(z) + \eta_v^2 q_{v,w}^{(2)}(z) + \frac{\eta_v}{\sqrt{nh}}q_{v,w}^{(3)}(z) \\ & - \mathbb{1}\{v \neq w\}\rho^{1+\hat{k}}(\Omega_1 + \rho^{\hat{k}}\Omega_2)\frac{\phi(z)}{2}z, \end{aligned}$$

we have

$$\sup_{z \in \mathbb{R}} |\mathbb{P}[T_{v,w} < z] - F_{v,w}(z)| = o\left((nh)^{-1} + (nh)^{-1/2}\eta_v + \eta_v^2 + \mathbb{1}\{v \neq w\}\rho^{1+2\hat{k}}\right).$$

To use this result to find the expansion of the error in coverage probability of the Normal-based confidence interval, the function  $F_{v,w}(z)$  is simply evaluated at the two endpoints of the interval. (Note: if the confidence interval were instead constructed with the bootstrap, a few additional steps are needed, but these do not alter any conclusions or results outside of constant terms.)

### S.I.6.1 Undersmoothing vs. Bias-Correction Exhausting all Smoothness

In general, we have assumed that the level of smoothness was large enough to be inconsequential in the analysis, and in particular this allowed for characterization of optimal bandwidth choices. In this section, in contrast, we take the level of smoothness to be binding, so that we can fully utilize the  $S$  derivatives *and* the Hölder condition to obtain the best possible rates of decay in coverage error for both undersmoothing and robust bias correction, but at the price of implementability: the leading bias constants can not be characterized, and hence feasible “optimal” bandwidths are not available.

For undersmoothing, the lowest bias is attained by setting  $\hat{k} > S$  (see Eqn. (10)), in which case the bias is only known to satisfy  $\mathbb{E}[\hat{f}] - f = O(h^{S+\varsigma})$  (i.e.,  $B_f$  is identically zero)



and bandwidth selection is not feasible. Note that this approach allows for  $\sqrt{nh}h^S \not\rightarrow 0$ , as  $\eta_{\text{us}} = O(\sqrt{nh}h^{S+\varsigma})$ .

Robust bias correction has several interesting features here. If  $k \leq S-2$  (the top two cases in Eqn. (11)), then the bias from approximating  $\mathbb{E}[\hat{f}] - f$  by  $B_f$ , that is not targeted by bias correction, dominates  $\eta_{\text{bc}}$  and prevents robust bias correction from performing as well as the best possible infeasible (i.e., oracle) undersmoothing approach. That is, even bias correction requires a sufficiently large choice of  $k$  in order to ensure the fastest possible rate of decay in coverage error: if  $k \geq S-1$ , robust bias correction can attain error decay rate as the best undersmoothing approach, and allow  $\sqrt{nh}h^S \not\rightarrow 0$ .

Within  $k \geq S-1$ , two cases emerge. On the one hand, if  $k = S-1$  or  $S$ , then  $B_f$  is nonzero and  $f^{(k)}$  must be consistently estimated to attain the best rate. Indeed, more is required. From Eqn. (11), we will need a bounded, positive  $\rho$  to equalize the bias terms. This (again) highlights the advantage of robust bias correction, as the classical procedure would enforce  $\rho \rightarrow 0$ , and thus underperform. On the other hand,  $\rho \rightarrow 0$  will be required if  $k > S$  because (from the final case of (11)) we require  $\rho^{k-S} = O(h^\varsigma)$  to attain the same rate as undersmoothing. Note that we can accommodate  $b \not\rightarrow 0$  (but bounded). Interestingly,  $B_f$  is identically zero and  $\hat{B}_f$  merely adds noise to the problem, but this noise is fully accounted for by the robust standard errors, and hence does not affect the rates of coverage error (though the constants of course change). The  $\hat{f}^{(k)}$  in  $\hat{B}_f$  is *inconsistent* ( $f^{(k)}$  does not exist), but the nonvanishing bias of  $\hat{f}^{(k)}$  is dominated by  $h^k$ .

This discussion is summarized by the following result:

**Corollary 7.** *Let the conditions of Theorem 3 hold.*

(a) *If  $k > S$ , then*

$$\mathbb{P}[f \in I_{\text{us}}] = 1 - \alpha + \frac{1}{nh} \frac{\phi(z_{\frac{\alpha}{2}})}{f} q_1(K) \{1 + o(1)\} + O(nh^{1+2S+2\varsigma} + h^{S+\varsigma}).$$

(b) *If  $k \geq S-1$ , then*

$$\begin{aligned} \mathbb{P}[f \in I_{\text{rbc}}] = 1 - \alpha + \frac{1}{nh} \frac{\phi(z_{\frac{\alpha}{2}})}{f} q_1(M) \{1 + o(1)\} \\ + O(nh(h^{S+\varsigma} \vee h^k b^{S-k+\varsigma} \mathbf{1}_{\{k \leq S\}})^2 + (h^{S+\varsigma} \vee h^k b^{S-k+\varsigma} \mathbf{1}_{\{k \leq S\}})). \end{aligned}$$

## S.I.6.2 Multivariate Densities and Derivative Estimation

We now briefly present state analogues of our results, both for distributional convergence and Edgeworth expansions, that cover multivariate data and derivative estimation. The conceptual

discussion and implications are similar to those in the main text, once adjusted notationally to the present setting, and are hence omitted.

For a nonnegative integral  $d$ -vector  $q$  we adopt the notation that: (i)  $[q] = q_1 + \dots + q_d$ , (ii)  $g^{(q)}(x) = \partial^{[q]}g(x)/(\partial^{q_1}x_1 \dots \partial^{q_d}x_d)$ , (iii)  $k! = q_1! \dots q_d!$ , and (iv)  $\sum_{[q]=Q}$  for some integer  $Q \geq 0$  denotes the sum over all indexes in the set  $\{q : [q] = Q\}$ .

The parameter of interest is  $f^{(q)}(x)$ , for  $x \in \mathbb{R}^d$  and  $[q] \leq S$ . The estimator is

$$\hat{f}^{(q)}(x) = \frac{1}{nh^{d+[q]}} \sum_{i=1}^n K^{(q)}(X_{h,i}).$$

Note that here, and below for bias correction, we use a constant, diagonal bandwidth matrix, e.g.  $h \times I_d$ . This is for simplicity and comparability, and could be relaxed at notational expense.

The bias, for a given kernel of order  $k \leq S - [q]$  (we restrict attention to the case where  $S$  is large enough), is

$$h^k \sum_{k:[k+q]=k} \mu_{K,k} f^{(q+k)}(x) + o(h^k),$$

exactly mirroring Eqn. (10), where now  $\mu_{K,k}$  represents a  $d$ -dimensional integral. Bias estimation is straightforward, relying on estimates  $\hat{f}^{(q+k)}(x)$ , for all  $[k] = k - [q]$ . The form of  $\hat{f}_2^{(q)}(x) = \hat{f}^{(q)}(x) - \hat{B}_{f^{(q)}}(x)$  is now given by

$$\hat{f}_2^{(q)}(x) = \frac{1}{nh^{d+[q]}} \sum_{i=1}^n M_{(q)}(X_{h,i}) \quad \text{where} \quad M_{(q)}(u) = K^{(q)}(u) - (\rho)^{d+[q]+k} \sum_{[k]=k} \mu_{K,k} L^{(q+k)}(u),$$

exactly analogous to Eqn. (9).

With these changes in notation out of the way, we can (re-)define the generic framework for both estimators exactly as above. Dropping the point of evaluation  $x$ , for  $v \in \{1, 2\}$ , define the estimator as

$$\hat{f}_v^{(q)} = \frac{1}{nh^{d+[q]}} \sum_{i=1}^n N_v(X_{h,i}), \quad \text{where} \quad N_1(u) = K^{(q)}(u) \text{ and } N_2(u) = M_{(q)}(u);$$

the variance

$$\sigma_v^2 := nh^{d+[q]} \mathbb{V}[\hat{f}_v^{(q)}] = \frac{1}{h^d} \{ \mathbb{E} [N_v(X_{h,i})^2] - \mathbb{E} [N_v(X_{h,i})]^2 \}$$

and its estimator as

$$\hat{\sigma}_v^2 = \frac{1}{h^d} \left\{ \frac{1}{n} \sum_{i=1}^n [N_v(X_{h,i})^2] - \left[ \frac{1}{n} \sum_{i=1}^n N_v(X_{h,i}) \right]^2 \right\};$$

and the  $t$ -statistics, for  $1 \leq w \leq v \leq 2$ , as,

$$T_{v,w} := \frac{\sqrt{nh^{d+2[q]}} \left( \hat{f}_v^{(q)} - f^{(q)} \right)}{\hat{\sigma}_w}.$$

As before,  $T_{\text{us}} = T_{1,1}$ ,  $T_{\text{bc}} = T_{2,1}$ , and  $T_{\text{rbc}} = T_{2,2}$ .

The scaled bias  $\eta_v$  has the same general definition as well: the bias of the numerator of the  $T_{v,w}$ . In this case, given by

$$\eta_v = \sqrt{nh^{d+2[q]}} \left( \mathbb{E} \left[ \hat{f}_v^{(q)} \right] - f^{(q)}(x) \right).$$

The asymptotic order of  $\eta_v$  for different settings can be obtained straightforwardly via the obvious multivariate extensions of Equation (11) and the corresponding conclusion of Lemma 1.

First-order convergence is now given by the following result. the proof of which is standard.

**Lemma 3.** *Suppose appropriate multivariate versions of Assumptions S.I.3.1 and S.I.3.2 hold,  $nh^{d+2[q]} \rightarrow \infty$ ,  $\eta_v \rightarrow 0$ , and if  $v = 2$ ,  $\rho \rightarrow 0 + \bar{\rho} \mathbb{1}\{v = w\}$ . Then  $T_{v,w} \rightarrow_d \mathcal{N}(0, 1)$ .*

For the Edgeworth expansion, redefine

$$\nu_{v,w}(j, k, p) = \frac{1}{h^{d+[q] \mathbb{1}\{j+pk=1\}}} \mathbb{E} \left[ (N_v(u_i) - \mathbb{E}[N_v(u_i)])^j (N_w(u_i)^p - \mathbb{E}[N_w(u_i)^p])^k \right],$$

where  $u_i = (x - X_i)/h$ . The polynomials  $p_{v,w}^{(k)}(z)$  and  $q_{v,w}^{(k)}(z)$  are as given above, but using multivariate moments. The analogue of Theorem 3 is given by the following result, which can be proven following the same steps as in Section S.I.7.

**Theorem 4.** *Suppose appropriate multivariate versions of Assumptions S.I.3.1, S.I.3.2, and S.I.3.3 hold,  $nh^{d+2[q]}/\log(n) \rightarrow \infty$ ,  $\eta_v \rightarrow 0$ , and if  $v = 2$ ,  $\rho \rightarrow 0 + \bar{\rho} \mathbb{1}\{v = w\}$ . Then for*

$$\begin{aligned} F_{v,w}(z) = & \Phi(z) + \frac{1}{\sqrt{nh^d}} p_{v,w}^{(1)}(z) + \sqrt{\frac{h^{d+2[q]}}{n}} p_{v,w}^{(2)}(z) + \eta_v p_{v,w}^{(3)}(z) + \frac{1}{nh^d} q_{v,w}^{(1)}(z) + \eta_v^2 q_{v,w}^{(2)}(z) + \frac{\eta_v}{\sqrt{nh^d}} q_{v,w}^{(3)}(z) \\ & + \mathbb{1}\{v \neq w\} \rho^{d+k+[q]} (\Omega_1 + \rho^{k+[q]} \Omega_2) \frac{\phi(z)}{2} z, \end{aligned}$$

we have

$$\sup_{z \in \mathbb{R}} |\mathbb{P}[T_{v,w} < z] - F_{v,w}(z)| = o\left(((nh^d)^{-1/2} + \eta_v)^2 + \mathbb{1}\{v \neq w\} \rho^{d+2(k+[q])}\right).$$

The same conclusions reached in the main text continue to hold for multivariate and/or derivative estimation, both in terms of comparing undersmoothing, bias correction, and robust bias correction, as well as for inference-optimal bandwidth choices. In particular, it is straightforward that the MSE optimal bandwidth in general has the rate  $n^{-1/(d+2k+2[q])}$ , whereas the coverage error optimal choice is of order  $n^{-1/(d+k+[q])}$ . Note that these two fit the same pattern as in the univariate, level case, with  $k + [q]$  in place of  $k$  and  $d$  in place of one. One intuitive reason for the similarity is that the number of derivatives in question does not impact that variance or higher order moment terms of the expansion, *once the scaling is accounted for*. That is, for all averages beyond the first, for example of the kernel squared,  $\sqrt{nh^d}$  can be thought of as the effective sample size, since that is the multiplier which stabilizes averages.

## S.I.7 Proof of Main Result

Throughout  $C$  shall be a generic constant that may take different values in different uses. If more than one constant is needed,  $C_1, C_2, \dots$ , will be used. It will cause no confusion (as the notations never occur in the same place), but in the course of proofs we will frequently write  $s = \sqrt{nh}$ , which overlaps with the order of the kernel  $L$ .

The first step is to write  $T_{v,w}$  as a smooth function of sums of i.i.d. random variables plus a remainder term that is shown to be of higher order. In addition to the notation above, define

$$\gamma_{v,p} = h^{-1} \mathbb{E}[N_v(X_{h,i})^p] \quad \text{and} \quad \Delta_{v,j} = \frac{1}{s} \sum_{i=1}^n \left\{ N_v(X_{h,i})^j - \mathbb{E}[N_v(X_{h,i})^j] \right\}.$$

With this notation  $\hat{f}_v - \mathbb{E}[\hat{f}_v] = s^{-1} \Delta_{v,1}$ ,  $\sigma_w^2 = \mathbb{E}[\Delta_{w,1}^2] = \gamma_{w,2} - h\gamma_{w,1}^2$  and

$$\hat{\sigma}_w^2 - \sigma_w^2 = s^{-1} \Delta_{w,2} - h2\gamma_{w,1}s^{-1} \Delta_{w,1} - hs^{-2} \Delta_{w,1}^2. \quad (12)$$

By a change of variables

$$\gamma_{v,p} = h^{-1} \int N_v(X_{h,i})^p f(X_i) dX_i = \int N_v(u)^p f(x - uh) du = O(1).$$

Further, by construction  $\mathbb{E}[\Delta_{w,j}] = 0$  and

$$\begin{aligned}\mathbb{V}[\Delta_{w,j}] &= h^{-1} \mathbb{E} \left[ N_v(X_{h,i})^{2j} \right] - h^{-1} \mathbb{E} \left[ N_v(X_{h,i})^j \right]^2 \\ &\leq h^{-1} \mathbb{E} \left[ N_v(X_{h,i})^{2j} \right] \\ &= \gamma_{v,2j} = O(1).\end{aligned}$$

Returning to Eqn. (12) and applying Markov's inequality, we find that  $hs^{-2}\Delta_{w,1}^2 = n^{-1}\Delta_{w,1}^2 = O_p(n^{-1})$  and  $\hat{\sigma}_w^2 - \sigma_w^2 = s^{-1}O_p(1) - hO(1)s^{-1}O_p(1) - hs^{-2}O_p(1) = O_p(s^{-1})$ , whence  $|\hat{\sigma}_w^2 - \sigma_w^2|^2 = O_p(s^{-2})$ . Using these results preceded by a Taylor expansion, we have

$$\begin{aligned}\left(\frac{\hat{\sigma}_w^2}{\sigma_w^2}\right)^{-1/2} &= \left(1 + \frac{\hat{\sigma}_w^2 - \sigma_w^2}{\sigma_w^2}\right)^{-1/2} = 1 - \frac{1}{2} \frac{\hat{\sigma}_w^2 - \sigma_w^2}{\sigma_w^2} + \frac{3}{8} \frac{(\hat{\sigma}_w^2 - \sigma_w^2)^2}{\sigma_w^4} + o_p((\hat{\sigma}_w^2 - \sigma_w^2)^2) \\ &= 1 - \frac{1}{2\sigma_w^2} (s^{-1}\Delta_{w,2} - h2\gamma_{w,1}s^{-1}\Delta_{w,1}) + O_p(n^{-1} + s^{-2}).\end{aligned}$$

Combining this result with the fact that

$$T_{v,w} = \frac{\Delta_{v,1} + \eta_v}{\hat{\sigma}_w} = \frac{\Delta_{v,1}}{\hat{\sigma}_w} + \frac{\eta_v}{\sigma_w} \left(\frac{\hat{\sigma}_w^2}{\sigma_w^2}\right)^{-1/2},$$

we have

$$\mathbb{P}[T_{v,w} < z] = \mathbb{P}\left[\tilde{T}_{v,w} - R_{v,w} < z - \frac{\eta_v}{\sigma_w}\right], \quad (13)$$

where

$$\tilde{T}_{v,w} = \frac{\Delta_{v,1}}{\hat{\sigma}_w} - \frac{\eta_v}{2\sigma_w^3} (s^{-1}\Delta_{w,2} - h2\gamma_{w,1}s^{-1}\Delta_{w,1})$$

and is a smooth function of sums of i.i.d. random variables and the remainder term is

$$R_{v,w} = \frac{\eta_v}{\sigma_w} \left( hs^{-2} \frac{\Delta_{w,1}^2}{2\sigma_w^2} + \frac{3}{8} \frac{(\hat{\sigma}_w^2 - \sigma_w^2)^2}{\sigma_w^4} + o_p((\hat{\sigma}_w^2 - \sigma_w^2)^2) \right).$$

Next we apply the delta method, see [Hall \(1992a, Chapter 2.7\)](#) or [Andrews \(2002, Lemma 5\(a\)\)](#). It will be true that

$$\mathbb{P}[T_{v,w} < z] = \mathbb{P}\left[\tilde{T}_{v,w} < z - \frac{\eta_v}{\sigma_w}\right] + o(s^{-2}) \quad (14)$$

if it can be shown that  $s^2\mathbb{P}[|R_{v,w}| > \varepsilon^2 s^{-2} \log(s)^{-1}] = o(1)$ .<sup>1</sup> This can be demonstrated by applying Bernstein's inequality to each piece of  $R_{v,w}$ , as the kernels  $K$  and  $L$ , and their derivatives, are bounded.

To apply this inequality to the first term of  $R_{v,w}$ , note that  $|N_w((x - X_i)/h)| \leq C_1$  and that  $\mathbb{V}[N_w((x - X_i)/h)] \leq C_2 h$ , for different constants, and so for  $\varepsilon > 0$  we have

$$\begin{aligned}
& s^2 \mathbb{P} \left[ \frac{\eta_v}{\sigma_w} h s^{-2} \frac{\Delta_{w,1}^2}{2\sigma_w^2} > \varepsilon^2 s^{-2} \log(s)^{-1} \right] \\
&= s^2 \mathbb{P} \left[ \left| \sum_{i=1}^n \{N_w(X_{h,i}) - \mathbb{E}[N_w(X_{h,i})]\} \right| > \varepsilon s^{-1} \log(s)^{-1/2} \left( \frac{2\sigma_w^3 n s^2}{\eta_v} \right)^{1/2} \right] \\
&= s^2 \mathbb{P} \left[ \left| \sum_{i=1}^n \{N_w(X_{h,i}) - \mathbb{E}[N_w(X_{h,i})]\} \right| > \varepsilon \left( \frac{2\sigma_w^3 n}{\eta_v \log(s)} \right)^{1/2} \right] \\
&\leq 2s^2 \exp \left\{ -\frac{1}{2} \frac{\varepsilon^2 2\sigma_w^3 n \eta_v^{-1} \log(s)^{-1}}{C_2 n h + \frac{1}{3} \varepsilon C_1 \sqrt{2\sigma_w^3 n / [\eta_v \log(s)]}} \right\} \\
&\leq s^2 \exp \left\{ -C \frac{\varepsilon^2 \log(s)^{-1}}{\eta h + \varepsilon \sqrt{\eta_v / [n \log(s)]}} \right\} \\
&\leq \exp \left\{ C_1 \log(s) \left[ 1 - C_2 \frac{\varepsilon^2}{\eta h \log(s)^2 + \varepsilon \sqrt{\eta_v \log(s)^3 / n}} \right] \right\},
\end{aligned}$$

which tends to zero because  $\eta_v \rightarrow 0$  as  $n \rightarrow \infty$  is assumed. To see why, note first that the second term of the denominator automatically vanishes, as  $\eta_v \rightarrow 0$  and  $\log(s)^3/n \rightarrow 0$ . Second, suppose  $\eta_v^2 \asymp n h^\omega$  (for example, if  $\eta_{\text{us}} \asymp s h^k$ , then  $\omega = 1 + 2k$ ) and the first term diverges, it must be that  $h$  is at least as large (in order) as

$$\left( \frac{1}{n \log(s)^4} \right)^{1/(2+\omega)},$$

which makes the requirement that  $\eta_v \rightarrow 0$  equivalent to

$$\eta_v^2 \asymp n h^\omega = n^{1-\omega/(2+\omega)} \log(s)^{-4\omega/(2+\omega)} \rightarrow 0,$$

which is impossible. The remaining terms of  $R_{v,w}$ , characterized using Eqn. (12), are handled in exactly the same way. This establishes Eqn. (14).

Next, the proofs of (Hall, 1992a, Chapters 4.4 and 5.5) show that  $\tilde{T}_{v,w}$  has an Edgeworth expansion valid through  $o(s^{-2} + s^{-1}\eta_v + \eta_v^2)$ . Thus, for a smooth function  $G(z)$  we can write

---

<sup>1</sup>Here,  $s^{-2} \log(s)^{-1}$  may be replaced with any sequence that is  $o(s^{-2} + \eta_v^2 + s^{-1}\eta_v)$ .

$\mathbb{P}[\tilde{T}_{v,w} < z] = G(z) + o(s^{-2} + s^{-1}\eta_v + \eta_v^2)$ . Therefore

$$\mathbb{P}\left[\tilde{T}_{v,w} < z - \frac{\eta_v}{\sigma_w}\right] = \mathbb{P}\left[\tilde{T}_{v,w} < z\right] - \frac{\eta_v}{\sigma_w}G^{(1)}(z) + o(s^{-2} + s^{-1}\eta_v + \eta_v^2). \quad (15)$$

The final result now follows by combining Equations (13), (14), and (15) with the terms of the expansion computed below.  $\square$

### S.I.7.1 Computing the Terms of the Expansion

Identifying the terms of the expansion is a matter of straightforward, if tedious, calculation. The first four cumulants of  $T_{v,w}$  must be calculated, which are functions of the first four moments. In what follows, we give a short summary. Note well that we always discard higher-order terms for brevity, and to save notation we will write  $\stackrel{o}{=}$  to stand in for “equal up to  $o((nh)^{-1} + (nh)^{-1/2}\eta_v + \eta_v^2 + \mathbb{1}\{v \neq w\}\rho^{1+2k})$ ”.

Referring to the Taylor expansion above, for the purpose of computing moments and cumulants, we can use

$$T_{v,w} \approx \left(\frac{\Delta_{v,1}}{\sigma_w} + \frac{\eta_v}{\sigma_w}\right) \left(1 - \frac{s^{-1}\Delta_{w,2}}{2\sigma_w} + \frac{h\gamma_{w,1}s^{-1}\Delta_{w,1}}{\sigma_w} + \frac{3}{8}\frac{s^{-2}\Delta_{w,2}^2}{\sigma_w^2}\right).$$

Moments of the two sides agree up to the requisite order. Straightforward moment calculations then give

$$\begin{aligned} \mathbb{E}[T_{v,w}] &\stackrel{o}{=} \frac{s^{-1}\mathbb{E}[\Delta_{v,1}\Delta_{w,2}]}{2\sigma_w^3} + \frac{hs^{-1}\gamma_{w,1}\mathbb{E}[\Delta_{v,1}\Delta_{w,1}]}{\sigma_w^3} + \frac{3s^{-2}\mathbb{E}[\Delta_{v,1}\Delta_{w,2}^2]}{8\sigma_w^5} + \frac{\eta_v}{\sigma_w} + \frac{3s^{-2}\eta_v\mathbb{E}[\Delta_{w,2}^2]}{8\sigma_w^5} \\ &\stackrel{o}{=} -s^{-1}\frac{\nu_{v,w}(1,1,2)}{2\sigma_w^3} + \frac{hs^{-1}\gamma_{w,1}\nu_{v,w}(1,1,1)}{\sigma_w^3} + \frac{\eta_v}{\sigma_w}, \\ \mathbb{E}[T_{v,w}^2] &\stackrel{o}{=} \frac{\mathbb{E}[\Delta_{v,1}^2]}{\sigma_w^2} + s^{-2}\frac{\mathbb{E}[\Delta_{v,1}^2\Delta_{w,2}^2]}{\sigma_w^6} + s^{-1}\frac{\mathbb{E}[\Delta_{v,1}^2\Delta_{w,2}]}{\sigma_w^4} + 2hs^{-1}\frac{\gamma_{w,1}\mathbb{E}[\Delta_{v,1}^2\Delta_{w,1}]}{\sigma_w^2} \\ &\quad - \eta_v s^{-1}\frac{2\mathbb{E}[\Delta_{v,1}\Delta_{w,2}]}{\sigma_w^4} + \eta_v hs^{-1}\frac{4\gamma_{w,1}\mathbb{E}[\Delta_{v,1}\Delta_{w,1}]}{\sigma_w^2} + \frac{\eta_v^2}{\sigma_w^2} \\ &\stackrel{o}{=} \frac{\sigma_v^2}{\sigma_w^2} + s^{-2}\frac{\sigma_v^2\nu_{v,w}(0,2,2)}{\sigma_w^6} + s^{-2}\frac{2\nu_{v,w}(1,1,2)^2}{\sigma_w^6} - s^{-2}\frac{\nu_{v,w}(2,1,2)^2}{\sigma_w^2} - \eta_v s^{-1}\frac{2\nu_{v,w}(1,1,2)}{\sigma_w^2} + \frac{\eta_v^2}{\sigma_w^2}, \\ \mathbb{E}[T_{v,w}^3] &\stackrel{o}{=} \frac{\mathbb{E}[\Delta_{v,1}^3]}{\sigma_w^3} - 3s^{-1}\frac{\mathbb{E}[\Delta_{v,1}^3\Delta_{w,2}]}{2\sigma_w^5} + 3hs^{-1}\frac{\gamma_{w,1}\mathbb{E}[\Delta_{v,1}^3\Delta_{w,1}]}{\sigma_w^5} + \eta_v\frac{3\mathbb{E}[\Delta_{v,1}^2]}{\sigma_w^3} - \eta_v s^{-1}\frac{9\mathbb{E}[\Delta_{v,1}^2\Delta_{w,2}]}{2\sigma_w^5} \end{aligned}$$

$$\stackrel{o}{=} s^{-1} \frac{\nu_v(3)}{\sigma_w^3} - s^{-1} \frac{9\nu_{v,w}(1,1,2)\sigma_v^2}{2\sigma_w^5} + hs^{-1} \frac{9\gamma_{w,1}\nu_{v,w}(1,1,1)}{\sigma_w^5} + \eta_v \frac{3\sigma_v^2}{\sigma_w^3},$$

and,

$$\begin{aligned} \mathbb{E}[T_{v,w}^4] &\stackrel{o}{=} \frac{\mathbb{E}[\Delta_{v,1}^4]}{\sigma_w^4} - s^{-1} \frac{2\mathbb{E}[\Delta_{v,1}^4 \Delta_{w,2}]}{\sigma_w^6} + 4hs^{-1} \frac{\gamma_{w,1}\mathbb{E}[\Delta_{v,1}^4 \Delta_{w,1}]}{\sigma_w^6} + s^{-2} \frac{3\mathbb{E}[\Delta_{v,1}^4 \Delta_{w,1}^2]}{\sigma_w^8} \\ &\quad + \eta_v \frac{4\mathbb{E}[\Delta_{v,1}^3]}{\sigma_w^4} - \eta_v s^{-1} \frac{8\mathbb{E}[\Delta_{v,1}^3 \Delta_{w,2}]}{\sigma_w^6} + \eta_v^2 \frac{6\mathbb{E}[\Delta_{v,1}^2]}{\sigma_w^4} \\ &\stackrel{o}{=} s^{-2} \frac{\nu_v(4)}{\sigma_w^4} + 3 \frac{\sigma_v^4}{\sigma_w^4} - s^{-2} \frac{8\nu_v(3)\nu_{v,w}(1,1,2) + 12\sigma_v^2\nu_{v,w}(2,1,2)}{\sigma_w^6} + s^{-2} \frac{9\sigma_v^4\nu_{v,w}(0,2,2)}{\sigma_w^8} \\ &\quad + s^{-2} \frac{36\sigma_v^2\nu_{v,w}(1,1,2)^2}{\sigma_w^8} + \eta_v s^{-1} \frac{4\nu_v(3)}{\sigma_w^4} - \eta_v s^{-1} \frac{24\sigma_v^2\nu_{v,w}(1,1,2)}{\sigma_w^6} + \eta_v^2 \frac{6\sigma_v^2}{\sigma_w^2}. \end{aligned}$$

The expansion now follows, formally, from the following steps. First, combining the above moments into cumulants. Second, these cumulants may be simplified using that

$$\frac{\sigma_v^2}{\sigma_w^2} = 1 + \mathbb{1}(w \neq v) (\rho^{1+k}\Omega_1 + \rho^{1+2k}\Omega_2)$$

and in all cases present

$$\nu_{v,w}(i, j, p) = f\vartheta_{N_v, i+jp} + o(1). \quad (16)$$

The second relation is readily proven for  $v = w$ , as  $\nu_{v,v}(i, j, p) = \mathbb{E}[N_v(X_{h,i})^{i+jp}] + O(h)$ , where the remainder represents products of expectations. In the case for  $v \neq w$ , we find  $\nu_{2,1}(i, j, p) = f\vartheta_{N_1, i+jp} + O(\rho^{1+k} + h)$ , and in this case  $\rho \rightarrow 0$  is assumed. For any term of a cumulant with a rate of  $(nh)^{-1}$ ,  $(nh)^{-1/2}\eta_v$ ,  $\eta_v^2$ , or  $\rho^{1+2k}$  (i.e., the extent of the expansion), these simplifications may be inserted as the remainder will be negligible. Note that this is exactly why the polynomials  $p_{v,w}^{(k)}$  do not simplify, while the  $q_{v,w}^{(k)}$  do. Third, with the cumulants in hand, the terms of the expansion are determined as described by e.g., [Hall \(1992a, Chapter 2\)](#).

Finally, for traditional bias correction, there are additional terms in the expansion (see discussion in the main text) representing the covariance of  $\hat{f}$  and  $\hat{B}_f$  (denoted by  $\Omega_1$ ) and the variance of  $\hat{B}_f$  ( $\Omega_2$ ). We now state their precise forms. These arise from the mismatch between the variance of the numerator of  $T_{bc}$  and the standardization used,  $\sigma_{us}^2$ , that is  $\sigma_{rbc}^2/\sigma_{us}^2$  is given by

$$\frac{nh\mathbb{V}[\hat{f} - \hat{B}_f]}{nh\mathbb{V}[\hat{f}]} = \frac{nh\mathbb{V}[\hat{f}] - 2nh\mathbb{C}[\hat{f}, \hat{B}_f] + nh\mathbb{V}[\hat{B}_f]}{nh\mathbb{V}[\hat{f}]} = 1 - 2\frac{nh\mathbb{C}[\hat{f}, \hat{B}_f]}{nh\mathbb{V}[\hat{f}]} + \frac{nh\mathbb{V}[\hat{B}_f]}{nh\mathbb{V}[\hat{f}]}.$$



This makes clear that  $\Omega_1$  and  $\Omega_2$  are the constant portions of the last two terms. First, for  $\Omega_1$ ,

$$\begin{aligned}\mathbb{C}[\hat{f}, \hat{B}_f] &= \mathbb{E} \left[ \left( \frac{1}{nh} \sum_{i=1}^n K(X_{h,i}) \right) \left( h^{\tilde{k}} \mu_{K,\tilde{k}} \frac{1}{nb^{1+\tilde{k}}} \sum_{i=1}^n L^{(\tilde{k})}(X_{b,i}) \right) \right] \\ &= h^{\tilde{k}} \mu_{K,\tilde{k}} \frac{1}{nb^{1+\tilde{k}}} \left\{ \mathbb{E} [h^{-1} K(X_{h,i}) L^{(\tilde{k})}(X_{b,i})] \right. \\ &\quad \left. - b \mathbb{E} [h^{-1} K(X_{h,i})] \mathbb{E} [b^{-1} L^{(\tilde{k})}(X_{b,i})] \right\} \\ &= \frac{\rho^{\tilde{k}} \mu_{K,\tilde{k}}}{nb} \left\{ \int f(x-uh) K(u) L^{(\tilde{k})}(u\rho) du - b \int f(x-uh) K(u) du \int f(x-ub) L^{(\tilde{k})}(u) du \right\}.\end{aligned}$$

Therefore

$$-2 \frac{nh \mathbb{C}[\hat{f}, \hat{B}_f]}{nh \mathbb{V}[\hat{f}]} = \rho^{1+\tilde{k}} \Omega_1,$$

where

$$\Omega_1 = -2 \frac{\mu_{K,\tilde{k}}}{\nu_1(2)} \left\{ \int f(x-uh) K(u) L^{(\tilde{k})}(u\rho) du - b \int f(x-uh) K(u) du \int f(x-ub) L^{(\tilde{k})}(u) du \right\}.$$

Note  $\nu_1(2) = \sigma_{\text{us}}^2$ . If we did not include  $\Omega_2$  in the Edgeworth expansion, i.e. we stopped at order  $\rho^{1+\tilde{k}}$ , then we could capture only the leading terms of  $\Omega_1$ , as follows, using that kernel integrates to 1 and  $\rho \rightarrow 0$ ,

$$\begin{aligned}\Omega_1 &= -2 \frac{\mu_{K,\tilde{k}}}{\nu_1(2)} \left\{ \int f(x-uh) K(u) L^{(\tilde{k})}(u\rho) du - b \int f(x-uh) K(u) du \int f(x-ub) L^{(\tilde{k})}(u) du \right\} \\ &= -2 \frac{\mu_{K,\tilde{k}}}{f(x) \vartheta_{K,2}^2 + O(h)} \left\{ f(x) L^{(\tilde{k})}(0) [1 + O(h + h\rho)] - b f(x)^2 \int L^{(\tilde{k})}(u) du [1 + O(b + h)] \right\} \\ &\rightarrow -2 \mu_{K,\tilde{k}} \vartheta_{K,2}^{-2} L^{(\tilde{k})}(0).\end{aligned}$$

Note that this matches the term [Hall \(1992b\)](#) calls  $w_2$ . We do not do this, for completeness. There are no other terms of up to order  $\rho^{1+2\tilde{k}}$ , so capturing the full contribution of  $\sigma_2^2/\sigma_1^2 - 1 = \sigma_{\text{rbc}}^2/\sigma_{\text{us}}^2 - 1$  is natural and informative.

Turning to  $\Omega_2$ , using the calculations in [Section S.I.4.1](#) (recall  $\tilde{k} = \tilde{k} \vee S$ ), we find that

$$\mathbb{V}[\hat{B}_f] = \frac{h^{2\tilde{k}}}{n} \mu_{K,\tilde{k}}^2 \left\{ \frac{1}{b^{1+2\tilde{k}}} \mathbb{E} [b^{-1} L^{(\tilde{k})}(X_{b,i})^2] - \left( \frac{1}{b^{1+\tilde{k}}} \mathbb{E} [L^{(\tilde{k})}(X_{b,i})] \right)^2 \right\}$$

$$= \frac{\rho^{2k} \mu_{K,k}^2}{nb} \left\{ \int f(x-ub) L^{(k)}(u)^2 du - b^{1+2\tilde{k}} \left( \int L^{(k-\tilde{k})}(u) f^{(\tilde{k})}(x-ub) du \right)^2 \right\},$$

and hence

$$\frac{nh\mathbb{V}[\hat{B}_f]}{nh\mathbb{V}[\hat{f}]} = \rho^{1+2k} \Omega_2 \quad \text{where} \quad \Omega_2 = \frac{\mu_{K,k}^2}{\nu_1(2)} \left\{ \int f(x-ub) L^{(k)}(u)^2 du - b^{1+2\tilde{k}} \left( \int L^{(k-\tilde{k})}(u) f^{(\tilde{k})}(x-ub) du \right)^2 \right\}.$$

The final piece will be  $b^{1+2S} f^{(k)}(x)^2 [1 + o(1)]$  if  $k \leq S$ . Substituting this is permitted because  $\rho^{1+2k}$  is the limit of the expansion, though it is not necessary to do, because this term is always higher order. Fully simplifying would yield

$$\Omega_2 = \mu_{K,k}^2 \vartheta_{K,2}^{-2} \vartheta_{L^{(k)},2},$$

which can be used in Theorem 3.

## S.I.8 Complete Simulation Results

To illustrate the gains from robust bias correction we conduct a Monte Carlo study to compare undersmoothing, traditional bias correction, and robust bias correction in terms coverage accuracy and interval length using several data-driven procedures to select the bandwidth. We generate  $n = 500$  observations from a true density  $f$  evaluated at  $x = \{-2, -1, 0, 1, 2\}$ . For the density, we consider:

Model 1 (Gaussian Density):  $x \sim \mathcal{N}(0, 1)$

Model 2 (Skewed Unimodal Density):  $x \sim \frac{1}{5}\mathcal{N}(0, 1) + \frac{1}{5}\mathcal{N}\left(\frac{1}{2}, \left(\frac{2}{3}\right)^2\right) + \frac{3}{5}\mathcal{N}\left(\frac{13}{12}, \left(\frac{5}{9}\right)^2\right)$

Model 3 (Bimodal Density):  $x \sim \frac{1}{2}\mathcal{N}\left(-1, \left(\frac{2}{3}\right)^2\right) + \frac{1}{2}\mathcal{N}\left(1, \left(\frac{2}{3}\right)^2\right)$

Model 4 (Asymmetric Bimodal Density):  $x \sim \frac{3}{4}\mathcal{N}(0, 1) + \frac{1}{4}\mathcal{N}\left(\frac{3}{2}, \left(\frac{1}{3}\right)^2\right)$

These models were previously analyzed in [Marron and Wand \(1992\)](#). They are plotted in Figure [S.I.1](#). In this simulation study we compare the performance of the confidence intervals defined by  $T_{\text{us}}$ ,  $T_{\text{bc}}$ , and  $T_{\text{rbc}}$ . For  $T_{\text{us}}$ , we take  $K$  to be the Epanechnikov kernel, while bias correction uses the Epanechnikov and MSE-optimal kernels for  $K$  and  $L^{(2)}$ , respectively. We consider two main data-drive bandwidth selectors. First, a Silverman rule-of-thumb alternative  $\hat{h}_{\text{rot}} = \hat{\sigma} 2.34 n^{-1/(2r+1)}$ . Second, the direct plug-in (DPI) for coverage error optimal

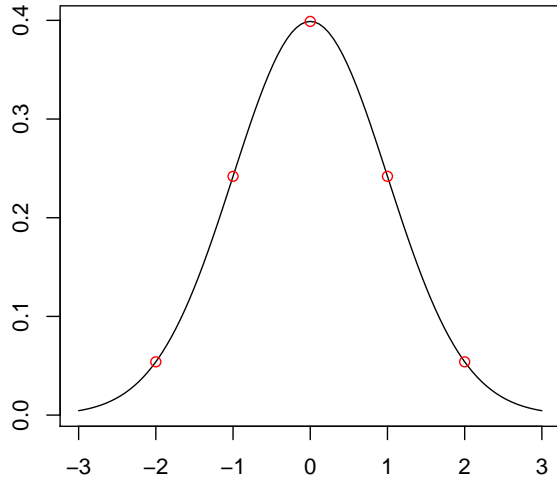
bandwidth  $\hat{h}_{\text{dpi}} = \hat{H}_{\text{dpi}} n^{-1/(r+3)}$ , where  $\hat{H}_{\text{dpi}}$  uses  $\hat{h}_{\text{rot}}$  as a pilot bandwidth to estimate  $f^{(r+2)}$  consistently. We also include the unfeasible, population value for  $h_{\text{mse}}$ .

Empirical coverage and length are reported in Tables S.I.2–S.I.5 (Panel A) using our two proposed data-driven bandwidth selectors, as well as the infeasible  $h_{\text{mse}}$ . The most obvious finding is that robust bias correction has accurate coverage for all bandwidth choices in all models. The intervals are generally longer than for undersmoothing, but neither undersmoothing nor traditional bias correction yield correct coverage outside of a few special cases (e.g., undersmoothing at the infeasible MSE-optimal bandwidth in Model 4). The DPI bandwidth selector generally results in slightly smaller bandwidths (on average). Summary statistics for the two fully data-driven bandwidths are shown in Panel B. The fact that the DPI bandwidth is slightly smaller is born out. It is also, in general, more variable.

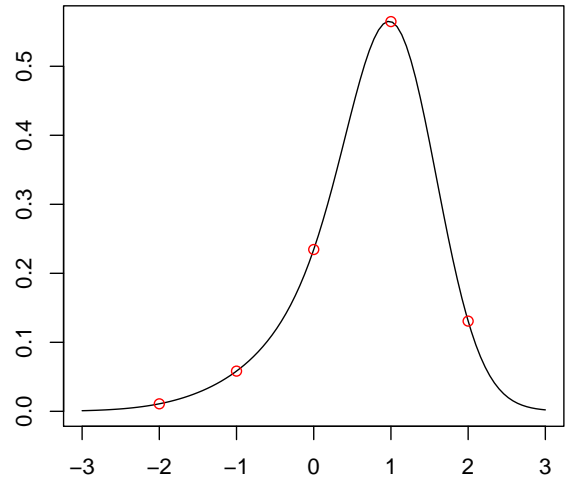
To illustrate the robustness to tuning parameter selection, Figures S.I.2–S.I.9 show coverage and length for all four models. The dotted vertical line shows the population MSE-optimal bandwidth for reference. These figures demonstrate the delicate balance required for undersmoothing to provide correct coverage, whereas for a wide range of bandwidths robust bias correction provides correct coverage. Further, interval length is not unduly inflated for bandwidths that provide correct coverage. Recall that robust bias correction can accommodate, and will optimally employ, a larger bandwidth, yielding higher precision. Further emphasizing the point of robustness, we depart from  $\rho = 1$  in Figures S.I.10 and S.I.11 to show coverage and length over a grid of  $h$  and  $\rho$ .

The simulation results for local polynomial regression reported in Section S.II.7 below bear out these same conclusions and study these issues in more detail, in particular interval length. All our methods are implemented in software available from the authors’ websites and via the R package `nprobust` available at <https://cran.r-project.org/package=nprobust>.

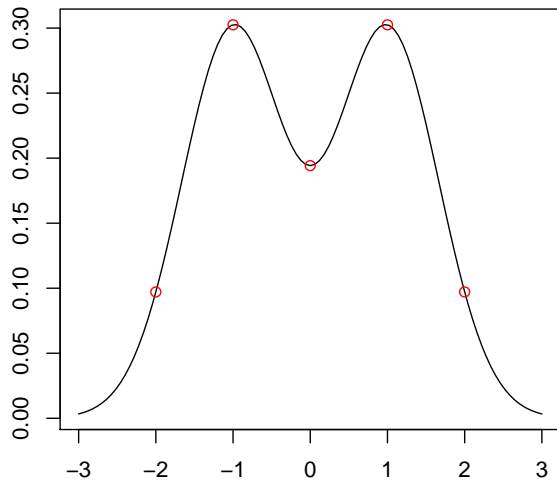
Figure S.I.1: Density Functions



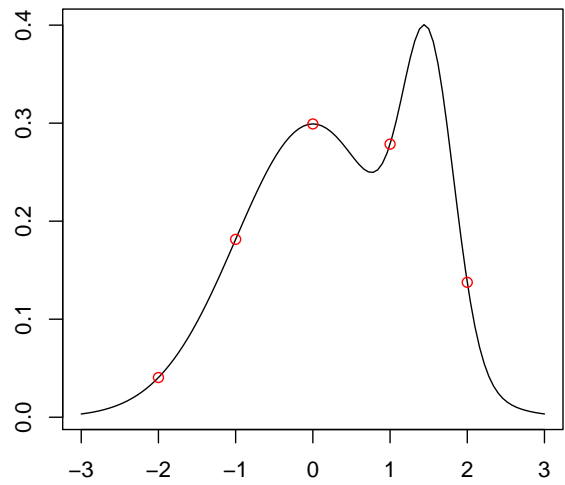
(a) Model 1



(b) Model 2



(c) Model 3



(d) Model 4

Table S.I.2: Simulations Results for Model 1

Panel A: Empirical Coverage and Average Interval Length of 95% Confidence Intervals

	Bandwidth	Empirical Coverage			Interval Length	
		US	BC	RBC	US	RBC
$x = -2$						
$h_{\text{mse}}$	0.677	91.5	81.1	93.9	0.039	0.055
$\hat{h}_{\text{rot}}$	0.674	90.6	81.2	93.9	0.039	0.055
$\hat{h}_{\text{dpi}}$	0.709	86.7	80.7	93.7	0.038	0.054
$x = -1$						
$h_{\text{mse}}$	0.677	94.5	78.0	94.3	0.069	0.109
$\hat{h}_{\text{rot}}$	0.674	94.5	78.0	94.0	0.069	0.109
$\hat{h}_{\text{dpi}}$	0.748	93.2	80.0	94.7	0.067	0.106
$x = 0$						
$h_{\text{mse}}$	0.677	85.5	74.8	95.0	0.078	0.132
$\hat{h}_{\text{rot}}$	0.674	84.4	74.8	94.9	0.078	0.132
$\hat{h}_{\text{dpi}}$	0.448	91.1	77.9	95.0	0.109	0.172
$x = 1$						
$h_{\text{mse}}$	0.677	94.9	78.3	94.5	0.069	0.109
$\hat{h}_{\text{rot}}$	0.674	94.7	78.3	94.6	0.069	0.109
$\hat{h}_{\text{dpi}}$	0.751	93.8	80.2	95.1	0.066	0.105
$x = 2$						
$h_{\text{mse}}$	0.677	92.0	83.2	94.3	0.038	0.055
$\hat{h}_{\text{rot}}$	0.674	91.2	83.2	94.3	0.039	0.055
$\hat{h}_{\text{dpi}}$	0.707	87.7	82.7	94.2	0.038	0.054

Panel B: Summary Statistics for the Estimated Bandwidths

	Pop. Par.	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Std. Dev.
$x = -2$								
$\hat{h}_{\text{rot}}$	0.677	0.5962	0.6597	0.6746	0.6744	0.6888	0.7503	0.021
$\hat{h}_{\text{dpi}}$	-	0.5597	0.6492	0.6856	0.7093	0.7354	2.394	0.11
$x = -1$								
$\hat{h}_{\text{rot}}$	0.677	0.5962	0.6597	0.6746	0.6744	0.6888	0.7503	0.021
$\hat{h}_{\text{dpi}}$	-	0.406	0.6196	0.7044	0.7484	0.8191	2.885	0.22
$x = 0$								
$\hat{h}_{\text{rot}}$	0.677	0.5962	0.6597	0.6746	0.6744	0.6888	0.7503	0.021
$\hat{h}_{\text{dpi}}$	-	0.3499	0.4084	0.4324	0.4478	0.4638	2.885	0.084
$x = 1$								
$\hat{h}_{\text{rot}}$	0.677	0.5962	0.6597	0.6746	0.6744	0.6888	0.7503	0.021
$\hat{h}_{\text{dpi}}$	-	0.4099	0.6241	0.7097	0.751	0.8227	2.885	0.21
$x = 2$								
$\hat{h}_{\text{rot}}$	0.677	0.5962	0.6597	0.6746	0.6744	0.6888	0.7503	0.021
$\hat{h}_{\text{dpi}}$	-	0.553	0.6506	0.6864	0.7071	0.7365	1.629	0.095

**Notes:**

(i) US = Undersmoothing, BC = Bias Corrected, RBC = Robust Bias Corrected.

(ii) Columns under “Bandwidth” report the average estimated bandwidths choices, as appropriate, for bandwidth  $h_n$ .

Table S.I.3: Simulations Results for Model 2

Panel A: Empirical Coverage and Average Interval Length of 95% Confidence Intervals

	Bandwidth	Empirical Coverage			Interval Length	
		US	BC	RBC	US	RBC
$x = -2$						
$h_{\text{mse}}$	0.454	90.5	79.1	85.5	0.021	0.028
$\hat{h}_{\text{rot}}$	0.551	91.7	79.3	87.6	0.019	0.026
$\hat{h}_{\text{dpi}}$	0.741	92.0	82.6	89.8	0.018	0.025
$x = -1$						
$h_{\text{mse}}$	0.454	94.4	80.4	92.9	0.048	0.069
$\hat{h}_{\text{rot}}$	0.551	94.1	80.7	93.1	0.043	0.063
$\hat{h}_{\text{dpi}}$	0.684	89.0	81.2	94.0	0.041	0.059
$x = 0$						
$h_{\text{mse}}$	0.454	94.9	81.0	95.1	0.089	0.135
$\hat{h}_{\text{rot}}$	0.551	93.0	80.9	95.0	0.079	0.122
$\hat{h}_{\text{dpi}}$	0.503	92.2	80.6	94.6	0.084	0.128
$x = 1$						
$h_{\text{mse}}$	0.454	83.8	76.1	94.9	0.115	0.193
$\hat{h}_{\text{rot}}$	0.551	62.2	74.1	94.4	0.097	0.169
$\hat{h}_{\text{dpi}}$	0.311	91.7	77.0	94.4	0.154	0.244
$x = 2$						
$h_{\text{mse}}$	0.454	90.6	82.1	93.8	0.071	0.104
$\hat{h}_{\text{rot}}$	0.551	79.4	81.8	94.4	0.064	0.095
$\hat{h}_{\text{dpi}}$	0.466	87.3	81.5	94.0	0.070	0.102

Panel B: Summary Statistics for the Estimated Bandwidths

	Pop. Par.	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Std. Dev.
$x = -2$								
$\hat{h}_{\text{rot}}$	0.454	0.477	0.5365	0.5504	0.5506	0.5649	0.6424	0.021
$\hat{h}_{\text{dpi}}$	-	0.411	0.5546	0.6643	0.741	0.8488	2.885	0.27
$x = -1$								
$\hat{h}_{\text{rot}}$	0.454	0.477	0.5365	0.5504	0.5506	0.5649	0.6424	0.021
$\hat{h}_{\text{dpi}}$	-	0.3724	0.5253	0.6442	0.6837	0.7681	2.885	0.23
$x = 0$								
$\hat{h}_{\text{rot}}$	0.454	0.477	0.5365	0.5504	0.5506	0.5649	0.6424	0.021
$\hat{h}_{\text{dpi}}$	-	0.3902	0.4693	0.4934	0.5027	0.5239	1.491	0.058
$x = 1$								
$\hat{h}_{\text{rot}}$	0.454	0.477	0.5365	0.5504	0.5506	0.5649	0.6424	0.021
$\hat{h}_{\text{dpi}}$	-	0.2545	0.2989	0.3093	0.3106	0.321	0.4302	0.017
$x = 2$								
$\hat{h}_{\text{rot}}$	0.454	0.477	0.5365	0.5504	0.5506	0.5649	0.6424	0.021
$\hat{h}_{\text{dpi}}$	-	0.3955	0.4474	0.4637	0.4661	0.4826	0.7111	0.028

**Notes:**

(i) US = Undersmoothing, BC = Bias Corrected, RBC = Robust Bias Corrected.

(ii) Columns under “Bandwidth” report the average estimated bandwidths choices, as appropriate, for bandwidth  $h_n$ .

Table S.I.4: Simulations Results for Model 3

Panel A: Empirical Coverage and Average Interval Length of 95% Confidence Intervals

	Bandwidth	Empirical Coverage			Interval Length	
		US	BC	RBC	US	RBC
$x = -2$						
$h_{\text{mse}}$	0.533	93.4	81.8	94.1	0.056	0.083
$\hat{h}_{\text{rot}}$	0.811	77.4	81.4	94.3	0.045	0.067
$\hat{h}_{\text{dpi}}$	0.839	72.1	78.0	92.7	0.045	0.066
$x = -1$						
$h_{\text{mse}}$	0.533	87.3	77.9	94.6	0.087	0.136
$\hat{h}_{\text{rot}}$	0.811	45.6	75.0	94.1	0.064	0.105
$\hat{h}_{\text{dpi}}$	0.467	88.9	78.7	94.4	0.095	0.147
$x = 0$						
$h_{\text{mse}}$	0.533	89.8	79.8	94.3	0.076	0.114
$\hat{h}_{\text{rot}}$	0.811	52.0	78.0	94.9	0.058	0.092
$\hat{h}_{\text{dpi}}$	0.646	79.6	78.2	94.5	0.067	0.103
$x = 1$						
$h_{\text{mse}}$	0.533	87.0	78.0	94.3	0.087	0.136
$\hat{h}_{\text{rot}}$	0.811	47.8	74.9	94.1	0.064	0.105
$\hat{h}_{\text{dpi}}$	0.467	88.9	78.6	94.3	0.095	0.147
$x = 2$						
$h_{\text{mse}}$	0.533	93.5	80.4	93.6	0.056	0.082
$\hat{h}_{\text{rot}}$	0.811	77.4	80.8	94.6	0.045	0.067
$\hat{h}_{\text{dpi}}$	0.839	72.7	77.4	92.3	0.045	0.066

Panel B: Summary Statistics for the Estimated Bandwidths

	Pop. Par.	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Std. Dev.
$x = -2$								
$\hat{h}_{\text{rot}}$	0.533	0.7394	0.7991	0.8112	0.8113	0.824	0.8851	0.019
$\hat{h}_{\text{dpi}}$	-	0.5429	0.7482	0.8023	0.839	0.8828	2.885	0.15
$x = -1$								
$\hat{h}_{\text{rot}}$	0.533	0.7394	0.7991	0.8112	0.8113	0.824	0.8851	0.019
$\hat{h}_{\text{dpi}}$	-	0.4027	0.4501	0.4646	0.4673	0.4813	0.6241	0.025
$x = 0$								
$\hat{h}_{\text{rot}}$	0.533	0.7394	0.7991	0.8112	0.8113	0.824	0.8851	0.019
$\hat{h}_{\text{dpi}}$	-	0.5623	0.6236	0.642	0.6465	0.6644	0.8991	0.034
$x = 1$								
$\hat{h}_{\text{rot}}$	0.533	0.7394	0.7991	0.8112	0.8113	0.824	0.8851	0.019
$\hat{h}_{\text{dpi}}$	-	0.4	0.4497	0.4643	0.467	0.4814	0.7607	0.025
$x = 2$								
$\hat{h}_{\text{rot}}$	0.533	0.7394	0.7991	0.8112	0.8113	0.824	0.8851	0.019
$\hat{h}_{\text{dpi}}$	-	0.6142	0.7495	0.8	0.8391	0.878	2.885	0.16

**Notes:**

(i) US = Undersmoothing, BC = Bias Corrected, RBC = Robust Bias Corrected.

(ii) Columns under “Bandwidth” report the average estimated bandwidths choices, as appropriate, for bandwidth  $h_n$ .

Table S.I.5: Simulations Results for Model 4

Panel A: Empirical Coverage and Average Interval Length of 95% Confidence Intervals

	Bandwidth	Empirical Coverage			Interval Length	
		US	BC	RBC	US	RBC
$x = -2$						
$h_{\text{mse}}$	0.395	94.3	81.4	92.5	0.043	0.062
$\hat{h}_{\text{rot}}$	0.739	90.3	82.2	93.9	0.032	0.046
$\hat{h}_{\text{dpi}}$	0.766	87.0	82.0	93.7	0.032	0.045
$x = -1$						
$h_{\text{mse}}$	0.395	94.4	80.1	94.1	0.086	0.128
$\hat{h}_{\text{rot}}$	0.739	94.6	79.2	94.2	0.059	0.091
$\hat{h}_{\text{dpi}}$	0.745	93.6	79.9	95.1	0.062	0.095
$x = 0$						
$h_{\text{mse}}$	0.395	94.1	79.4	94.8	0.105	0.161
$\hat{h}_{\text{rot}}$	0.739	85.2	77.0	95.0	0.069	0.112
$\hat{h}_{\text{dpi}}$	0.785	66.3	76.5	95.1	0.069	0.112
$x = 1$						
$h_{\text{mse}}$	0.395	93.2	79.6	95.3	0.104	0.158
$\hat{h}_{\text{rot}}$	0.739	67.0	73.2	93.6	0.068	0.112
$\hat{h}_{\text{dpi}}$	0.590	89.9	82.7	96.9	0.083	0.131
$x = 2$						
$h_{\text{mse}}$	0.395	90.7	82.4	94.2	0.079	0.115
$\hat{h}_{\text{rot}}$	0.739	33.8	74.0	91.9	0.057	0.085
$\hat{h}_{\text{dpi}}$	0.762	33.5	75.3	92.0	0.058	0.087

Panel B: Summary Statistics for the Estimated Bandwidths

	Pop. Par.	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Std. Dev.
$x = -2$								
$\hat{h}_{\text{rot}}$	0.395	0.6624	0.7261	0.7396	0.7395	0.7526	0.8147	0.02
$\hat{h}_{\text{dpi}}$	-	0.5778	0.7103	0.7466	0.7664	0.7949	2.885	0.11
$x = -1$								
$\hat{h}_{\text{rot}}$	0.395	0.6624	0.7261	0.7396	0.7395	0.7526	0.8147	0.02
$\hat{h}_{\text{dpi}}$	-	0.4469	0.5654	0.6646	0.7448	0.8531	2.885	0.26
$x = 0$								
$\hat{h}_{\text{rot}}$	0.395	0.6624	0.7261	0.7396	0.7395	0.7526	0.8147	0.02
$\hat{h}_{\text{dpi}}$	-	0.414	0.6018	0.7431	0.7855	0.8893	2.885	0.26
$x = 1$								
$\hat{h}_{\text{rot}}$	0.395	0.6624	0.7261	0.7396	0.7395	0.7526	0.8147	0.02
$\hat{h}_{\text{dpi}}$	-	0.4047	0.4915	0.5302	0.5896	0.6039	2.636	0.18
$x = 2$								
$\hat{h}_{\text{rot}}$	0.395	0.6624	0.7261	0.7396	0.7395	0.7526	0.8147	0.02
$\hat{h}_{\text{dpi}}$	-	0.4532	0.5822	0.6896	0.7617	0.8656	2.885	0.25

**Notes:**

(i) US = Undersmoothing, BC = Bias Corrected, RBC = Robust Bias Corrected.

(ii) Columns under “Bandwidth” report the average estimated bandwidths choices, as appropriate, for bandwidth  $h_n$ .



Figure S.I.2: Empirical Coverage of 95% Confidence Intervals - Model 1

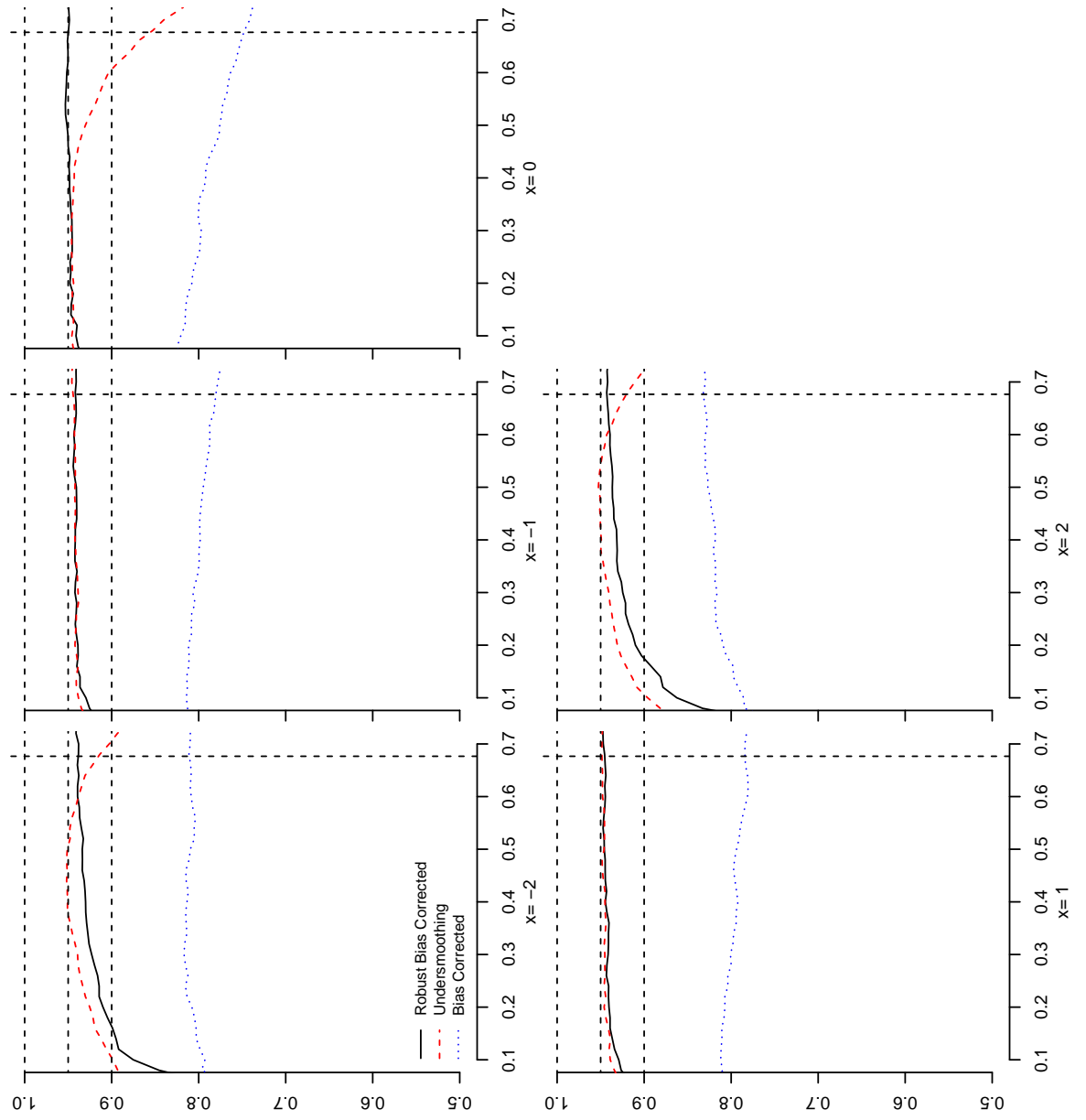


Figure S.I.3: Empirical Coverage of 95% Confidence Intervals - Model 2

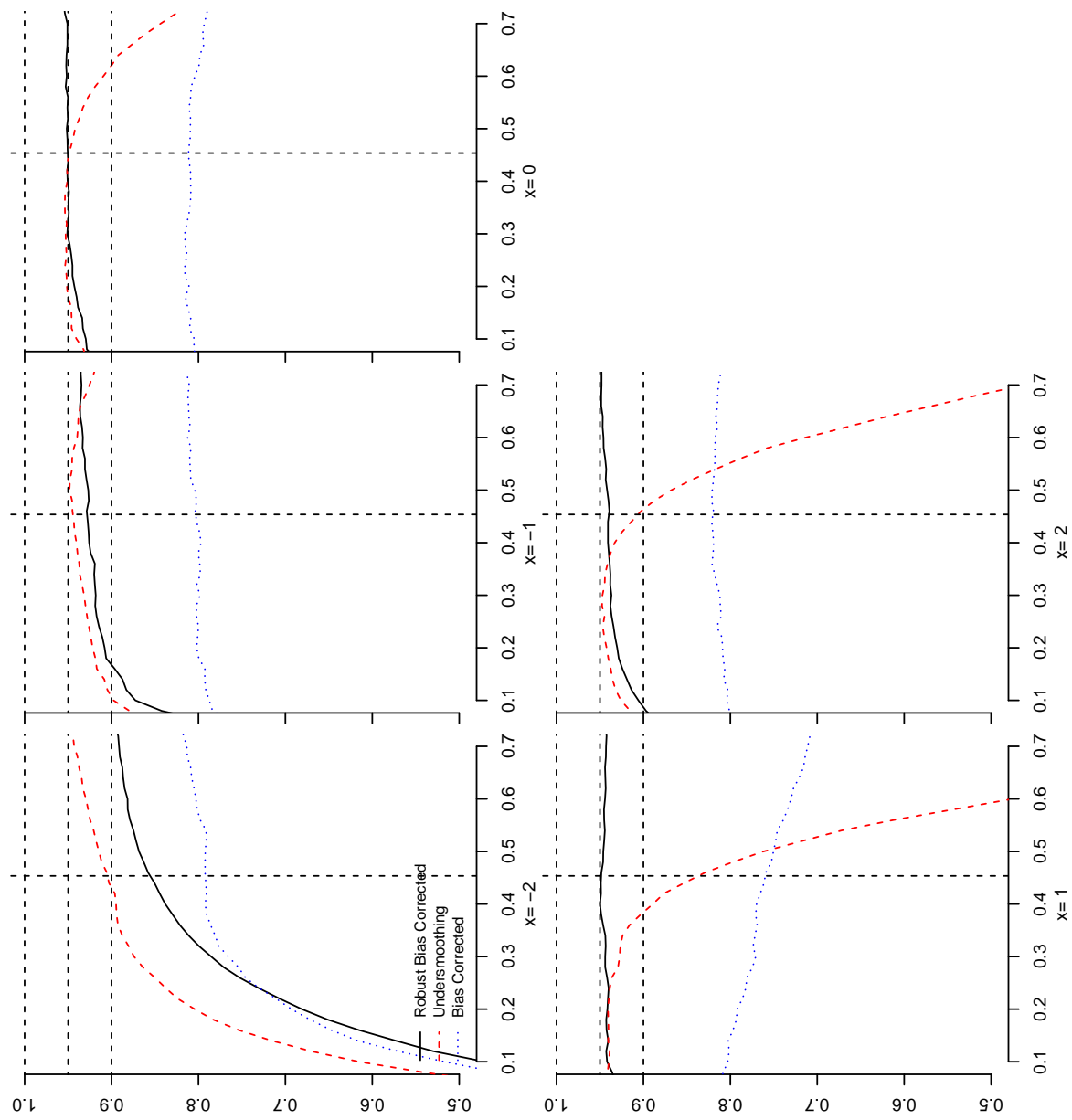


Figure S.I.4: Empirical Coverage of 95% Confidence Intervals - Model 3

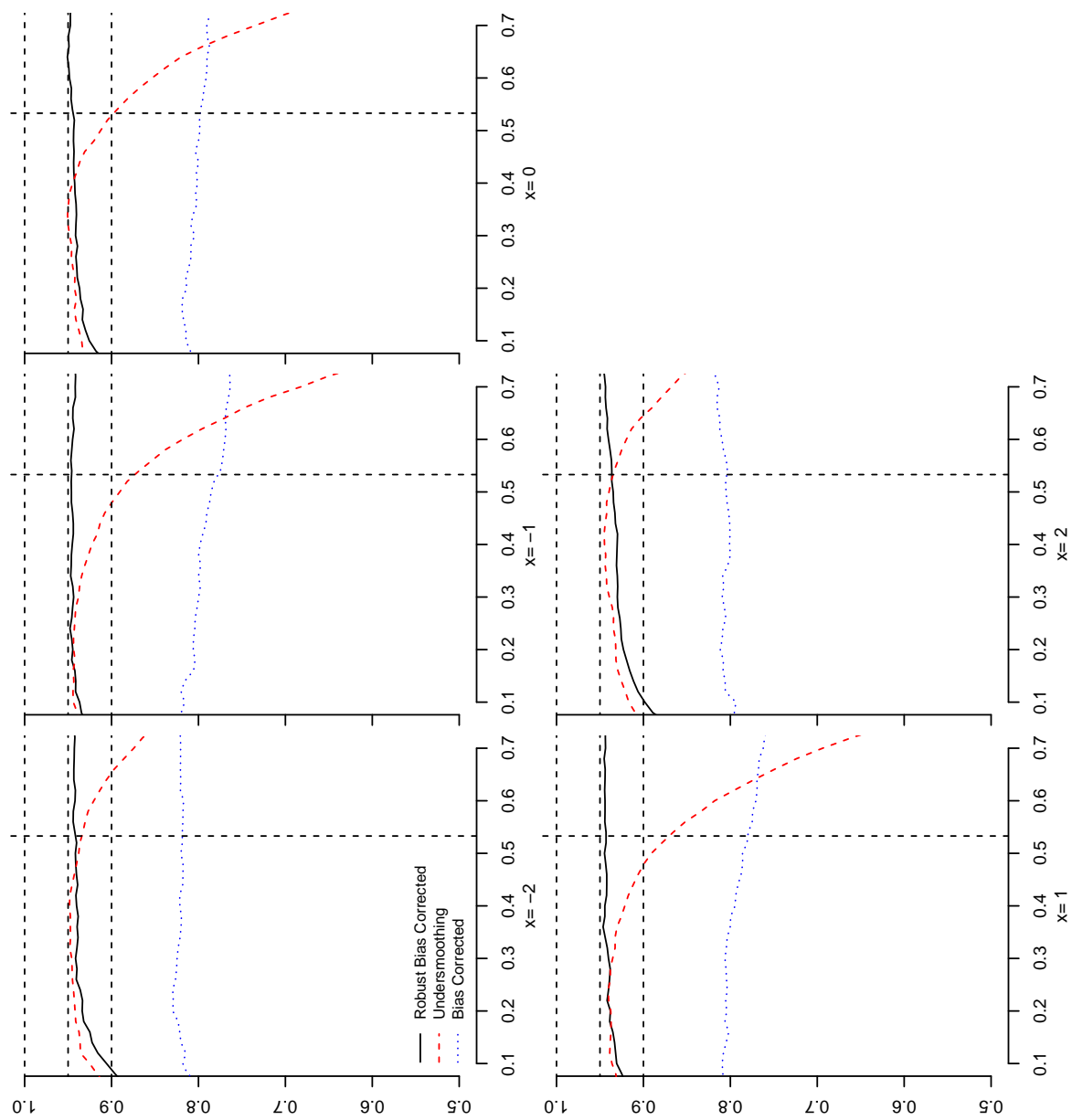


Figure S.I.5: Empirical Coverage of 95% Confidence Intervals - Model 4

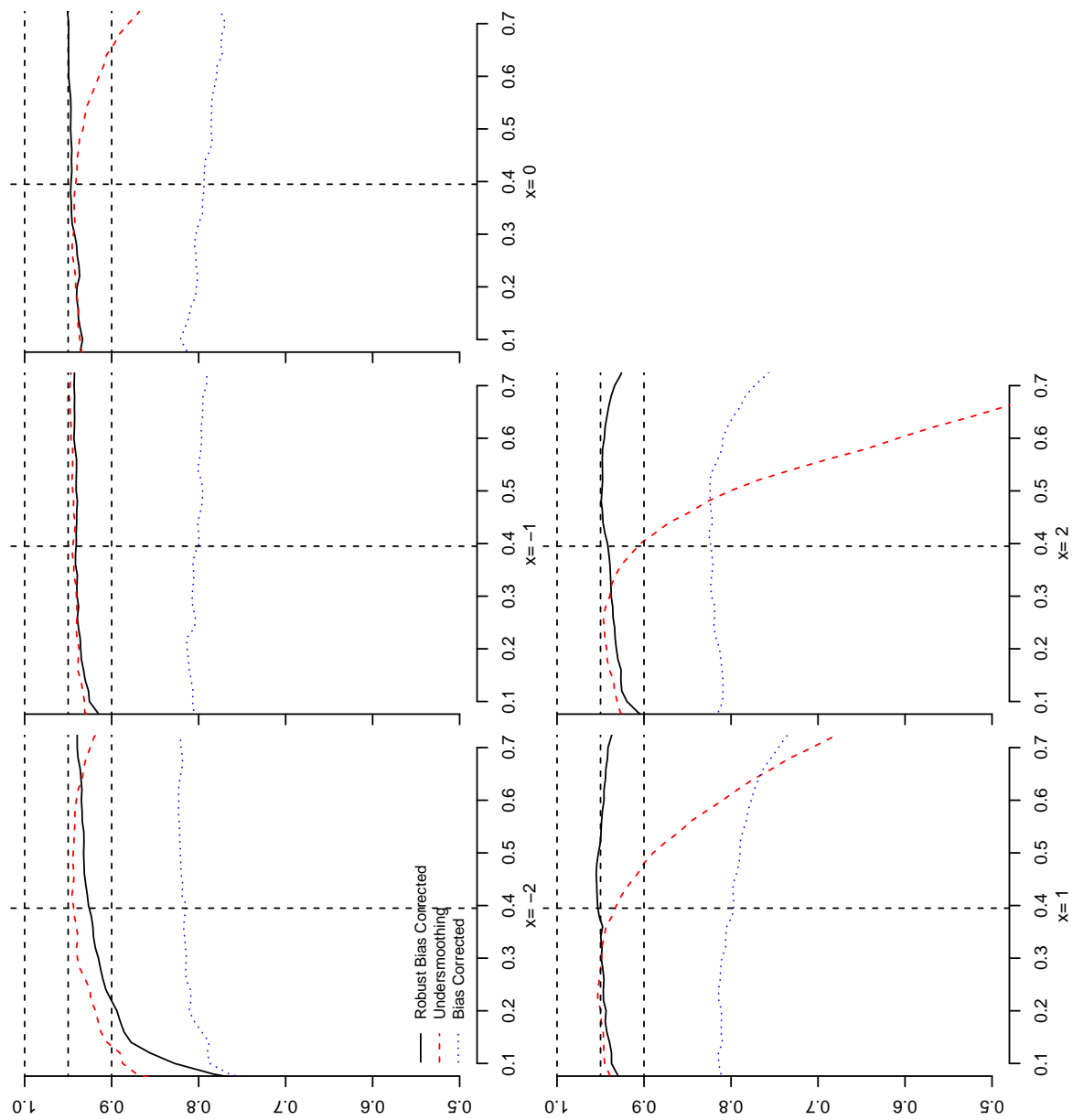


Figure S.I.6: Average Interval Length of 95% Confidence Intervals - Model 1

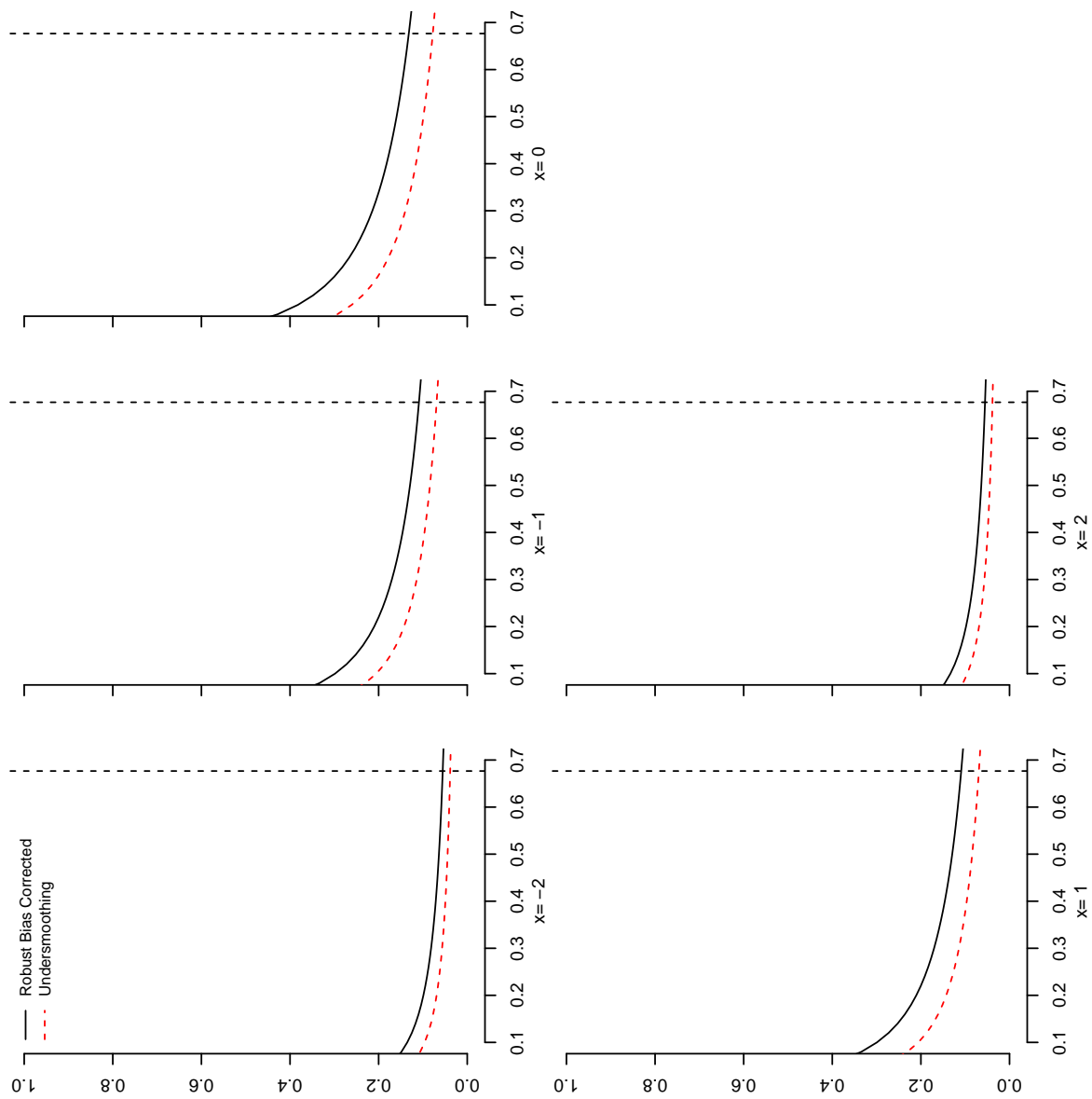


Figure S.I.7: Average Interval Length of 95% Confidence Intervals - Model 2

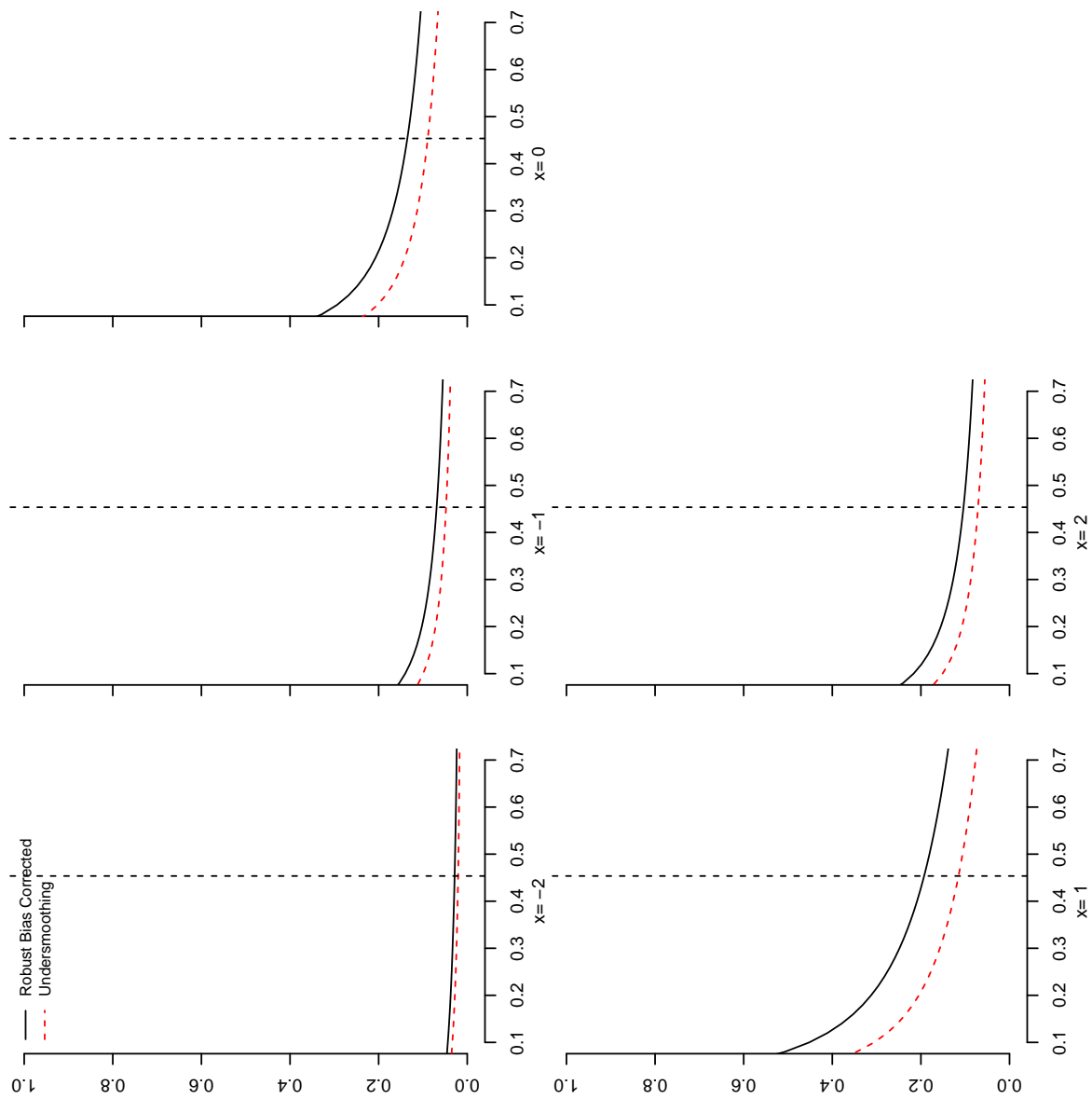


Figure S.I.8: Average Interval Length of 95% Confidence Intervals - Model 3

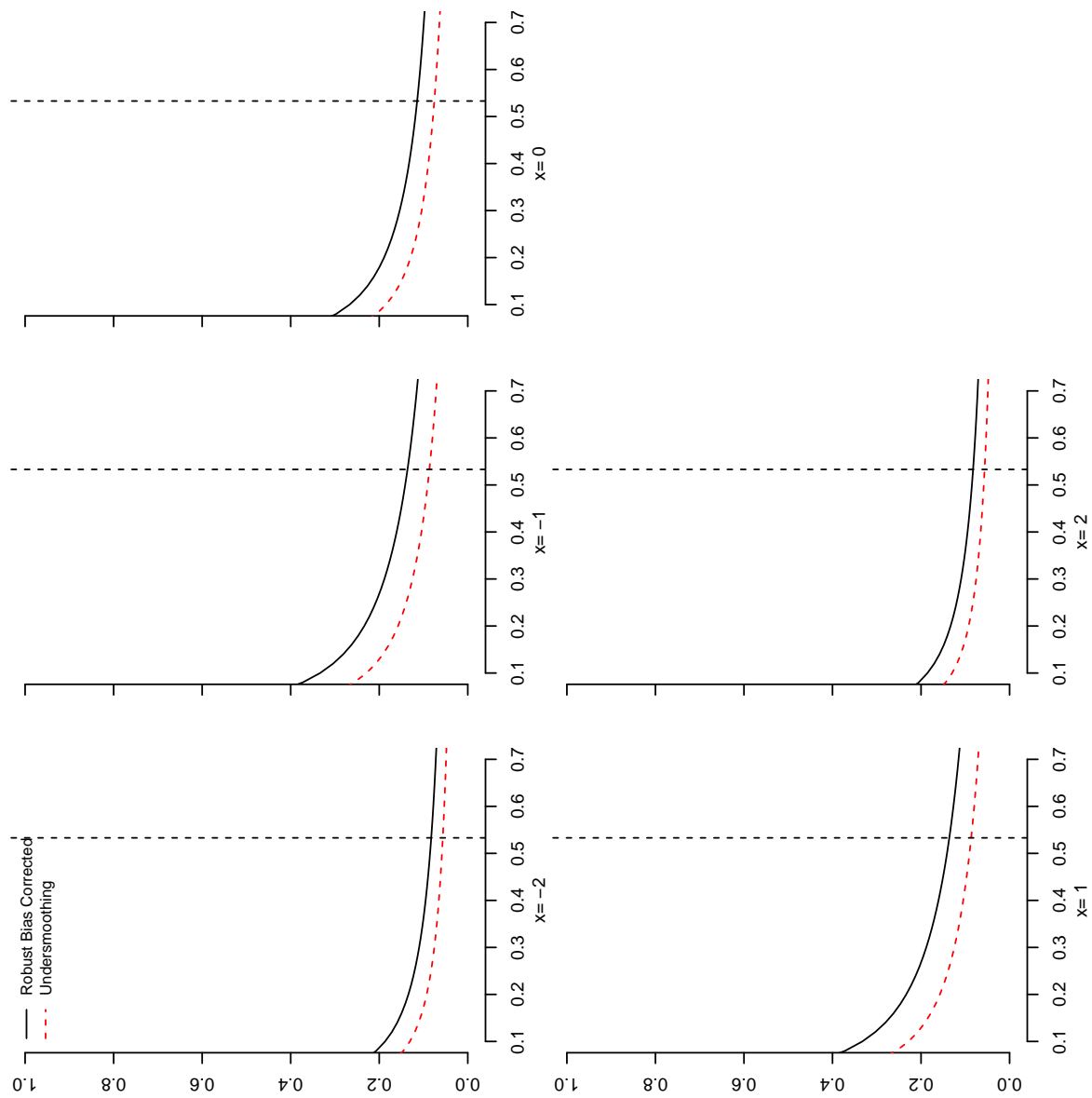


Figure S.I.9: Average Interval Length of 95% Confidence Intervals - Model 4

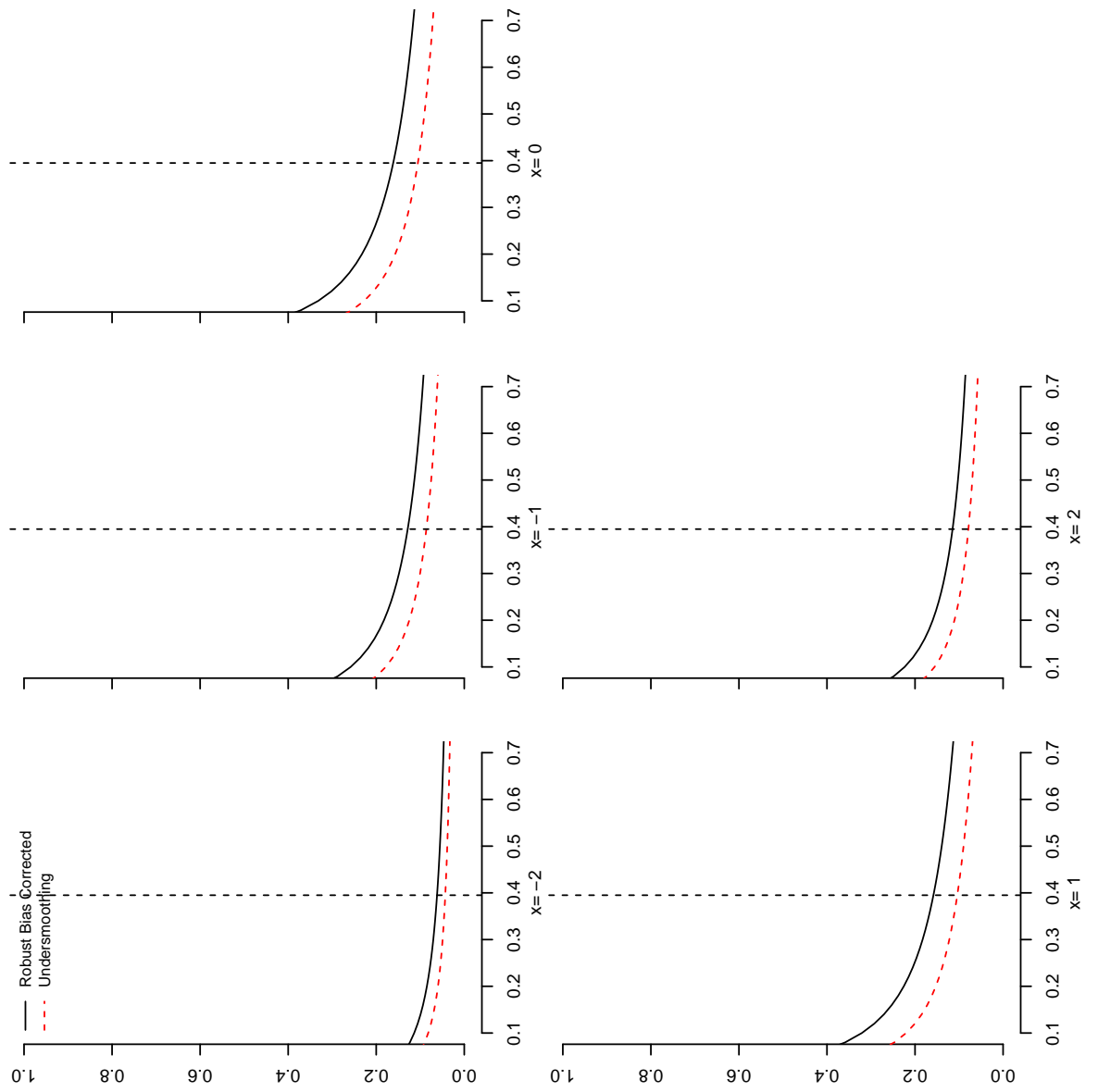
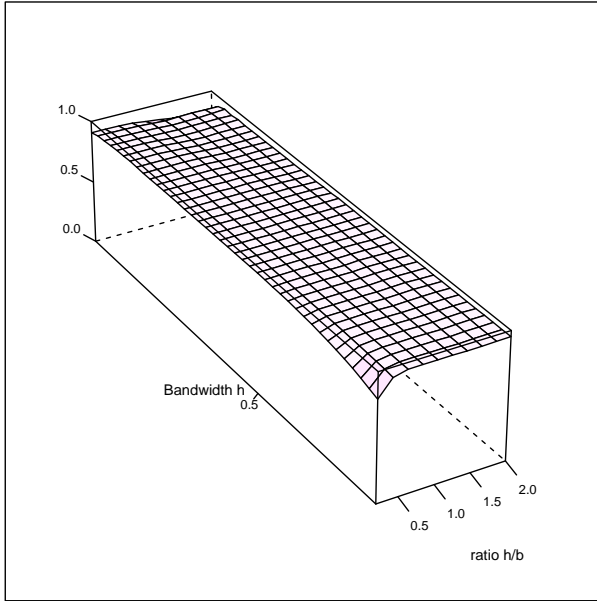
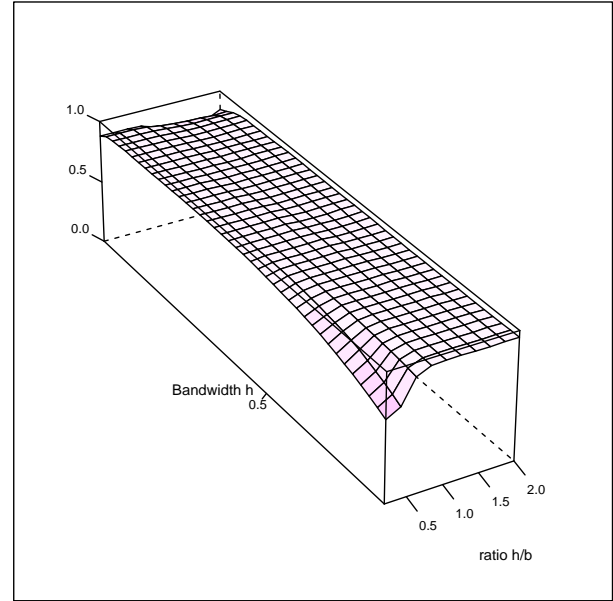




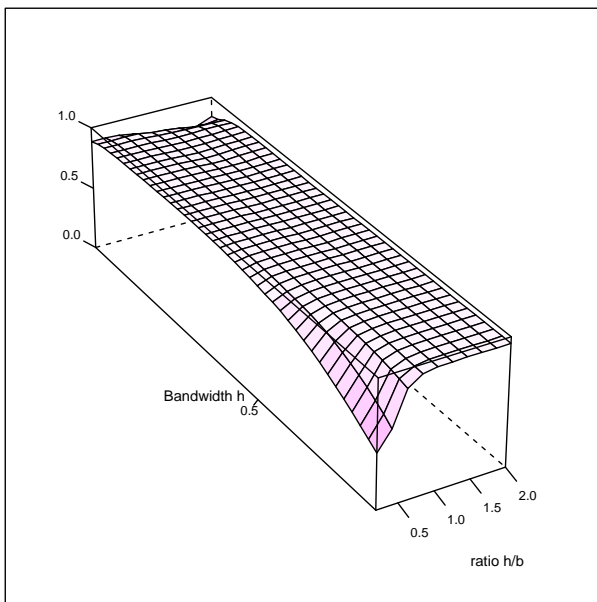
Figure S.I.10: Empirical Coverage of 95% Confidence Intervals ( $x = 0$ )



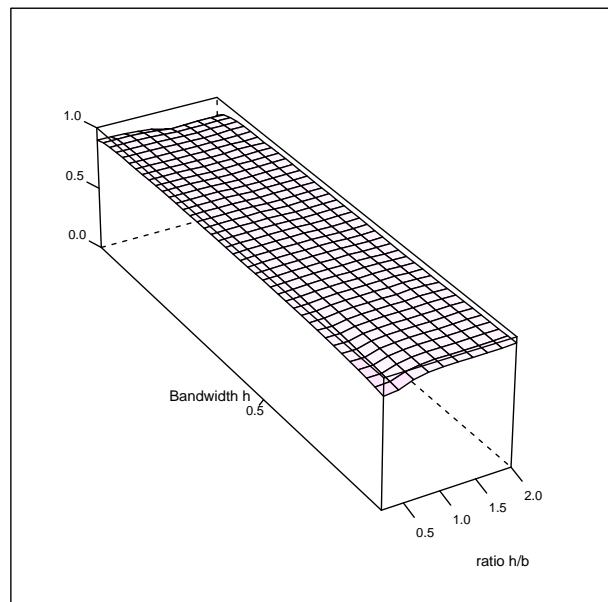
(a) Model 1



(b) Model 2

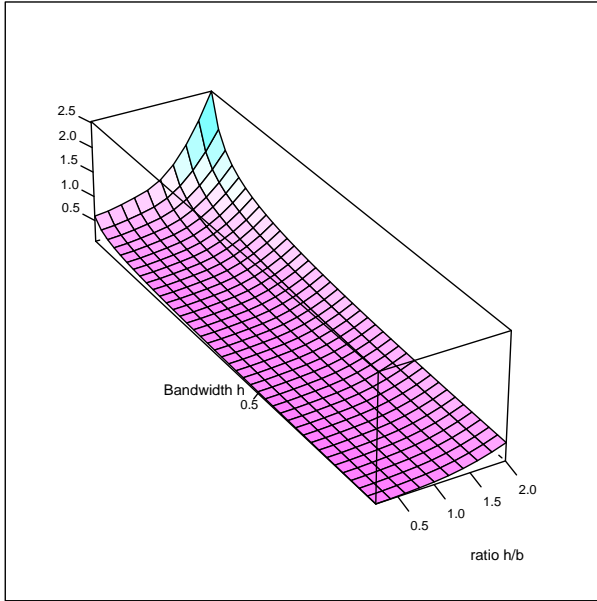


(c) Model 3

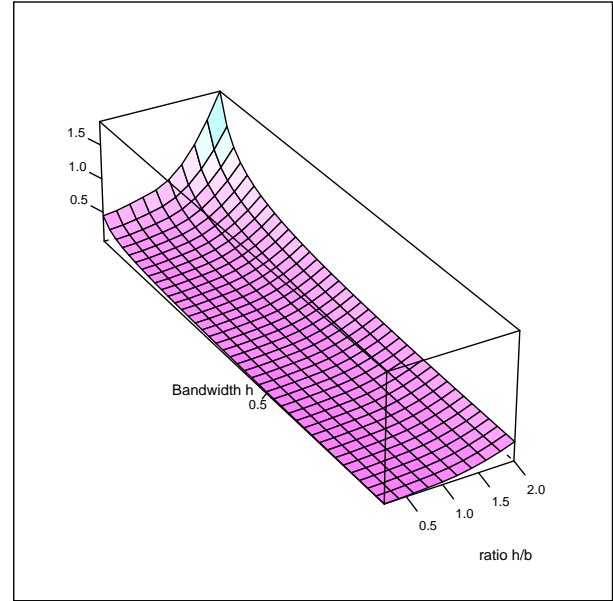


(d) Model 4

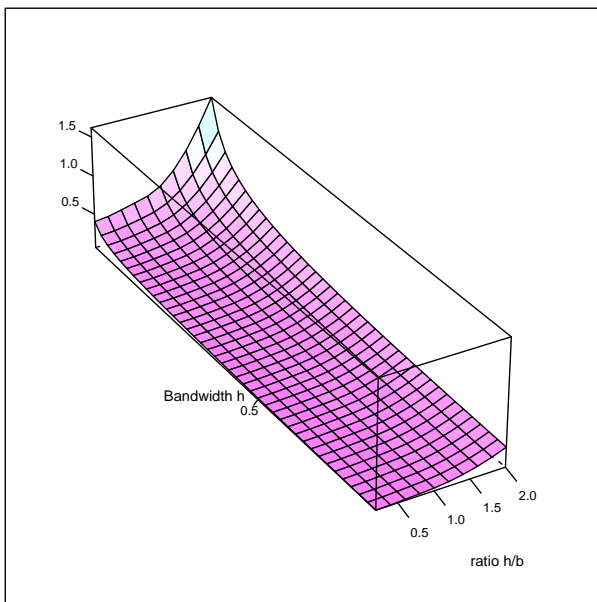
Figure S.I.11: Average Interval Length of 95% Confidence Intervals ( $x = 0$ )



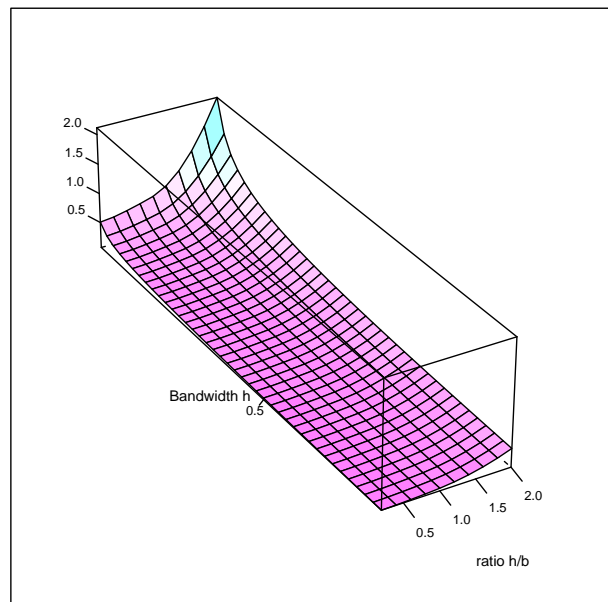
(a) Model 1



(b) Model 2



(c) Model 3



(d) Model 4

## Part S.II

# Local Polynomial Estimation and Inference

### S.II.1 Notation

Local polynomial regression is notationally demanding, and the Edgeworth expansions will be substantially more so. For ease of reference, we collect all notation here regardless of where it is introduced and used. Much of the notation is fully restated later, when needed. As such, this subsection is designed more for reference, and is not easily readable.

Throughout, a subscript  $p$  will generally refer to a quantity used to estimate  $m(x) = \mathbb{E}[Y_i|X_i = x]$ , while a subscript  $q$  will refer to the bias correction portion (the vectors  $e_0$  and  $e_{p+1}$  below are notable exceptions to this rule). Recall that  $p \geq 1$  is odd and  $q > p$  may be even or odd.

Throughout this section let  $X_{h,i} = (X_i - x)/h$  and similarly for  $X_{b,i}$ . The evaluation point is implicit here.

To save notation, products of functions will be written together, with only one argument. For example

$$(Kr_p r'_p)(X_{h,i}) := K(X_{h,i})r_p(X_{h,i})r'_p(X_{h,i})' = K\left(\frac{X_i - x}{h}\right)r_p\left(\frac{X_i - x}{h}\right)r'_p\left(\frac{X_i - x}{h}\right)',$$

and similarly for  $(Kr_p)(X_{h,i})$ ,  $(Lr_q)(X_{b,i})$ , etc.

All expectations are fixed- $n$  calculations. To give concrete examples of this notation ( $\Lambda_{p,k}$ ,  $R_p$ , and  $W_p$  are redefined below):

$$\Lambda_{p,k} = R'_p W_p [((X_1 - x)/h)^{p+k}, \dots, ((X_n - x)/h)^{p+k}]'/n = \frac{1}{nh} \sum_{i=1}^n (Kr_p)(X_{h,i}) X_{h,i}^{p+1}$$

and

$$\tilde{\Lambda}_{p,k} = \mathbb{E}[\Lambda_{p,k}] = h^{-1} \mathbb{E}[(Kr_p)(X_{h,j}) X_{h,i}^{p+k}] = h^{-1} \int_{\text{supp}\{X\}} K\left(\frac{X_i - x}{h}\right) r_p\left(\frac{X_i - x}{h}\right) \left(\frac{X_i - x}{h}\right)^{p+k} f(X_i) dX_i$$

Here the range of integration is explicit, but in general it will not be. This is important for boundary issues, where the notation is generally unchanged, and it is to be understood

that moments and moments of the kernel be replaced by the appropriate truncated version. Continuing this example, if  $\text{supp}\{X\} = [0, \infty)$  and  $x = 0$ , then by a change of variables

$$\tilde{\Lambda}_{p,k} = h^{-1} \int_{\text{supp}\{X\}} (Kr_p)(X_{h,j}) X_{h,i}^{p+k} f(X_i) dX_i = \int_0^\infty (Kr_p)(u) u^{p+k} f(-uh) du,$$

whereas if  $\text{supp}\{X\} = (-\infty, 0]$  and  $x = 0$ , then

$$\tilde{\Lambda}_{p,k} = \int_{-\infty}^0 (Kr_p)(u) u^{p+k} f(-uh) du.$$

For the remainder of this section, the notation is left generic.

For the proofs (Section [S.II.6](#)) we will frequently abbreviate  $s = \sqrt{nh}$ .

### S.II.1.1 Estimators, Variances, and Studentized Statistics

To define the estimator  $\hat{m}$  of  $m$  and the bias correction, begin by defining:

$$\begin{aligned} r_p(u) &= (1, u, u^2, \dots, u^p)', & R_p &= [r_p(X_{h,1}), \dots, r_p(X_{h,n})]', \\ W_p &= \text{diag}(h^{-1}K(X_{h,i}) : i = 1, \dots, n), & H_p &= \text{diag}(1, h^{-1}, h^{-2}, \dots, h^{-p}), \\ \Gamma_p &= R_p' W_p R_p / n, & \text{and} & \quad \Lambda_{p,k} = R_p' W_p [X_{h,1}^{p+k}, \dots, X_{h,n}^{p+k}]' / n, \end{aligned} \quad (17)$$

where  $\text{diag}(a_i : i = 1, \dots, n)$  denote the  $n \times n$  diagonal matrix constructed using the elements  $a_1, a_2, \dots, a_n$ . Note that in the main text  $\Lambda_{p,1}$  is denoted by  $\Lambda_p$ .

Similarly, define

$$\begin{aligned} r_q(u) &= (1, u, u^2, \dots, u^q)', & R_q &= [r_q(X_{b,1}), \dots, r_q(X_{b,n})]', \\ W_q &= \text{diag}(b^{-1}L(X_{b,i}) : i = 1, \dots, n), & H_q &= \text{diag}(1, b^{-1}, b^{-2}, \dots, b^{-q}), \\ \Gamma_q &= R_q' W_q R_q / n, & \text{and} & \quad \Lambda_{q,k} = R_q' W_q [X_{b,1}^{q+k}, \dots, X_{b,n}^{q+k}]' / n, \end{aligned} \quad (18)$$

These are identical, but substituting  $q$ ,  $L$ , and  $b$  in place of  $p$ ,  $K$ , and  $h$ , respectively. Note that some dimensions change but other do not: for example,  $W_p$  and  $W_q$  are both  $n \times n$ , but  $\Gamma_p$  is  $(p+1)$  square whereas  $\Gamma_q$  is  $(q+1)$ .

Denote by  $e_0$  the  $(p+1)$ -vector with a one in the first position and zeros in the remaining and  $Y = (Y_1, \dots, Y_n)'$ . The local polynomial estimator of  $m(x) = \mathbb{E}[Y_i | X_i = x]$  is

$$\hat{m} = e_0' \hat{\beta}_p = e_0' H_p \Gamma_p^{-1} R_p' W_p Y / n,$$

where

$$\hat{\beta}_p = \arg \min_{b \in \mathbb{R}^{p+1}} \frac{1}{nh} \sum_{i=1}^n (Y_i - r_p(X_i - x)'b)^2 K(X_{h,i}) = H_p \Gamma_p^{-1} R_p' W_p Y / n.$$

If we define  $\check{R} = [r_p(X_1 - x), \dots, r_p(X_n - x)]'$  and  $M = [m(X_1), \dots, m(X_n)]'$ , then we can split  $\hat{m} - m$  into the variance and bias terms

$$\hat{m} - m = e_0' \Gamma_p^{-1} R_p' W_p (Y - M) / n + e_0' \Gamma_p^{-1} R_p' W_p (M - \check{R} \beta_p) / n.$$

This will be useful in the course of the proofs.

The conditional bias is given by

$$\mathbb{E}[\hat{m}|X_1, \dots, X_n] - m = h^{p+1} m^{(p+1)} \frac{1}{(p+1)!} e_0' \Gamma_p^{-1} \Lambda_{p,1} + o_P(h^{p+1}). \quad (19)$$

(Recall that in the main paper,  $\Lambda_{p,1}$  is denoted  $\Lambda_p$ .) This result is valid for  $p$  odd, our main focus, but also for  $p$  even at boundary points.

Denote by  $e_{p+1}$  the  $(q+1)$ -vector with one in the  $p+2$  position, and zeros in the rest. Then we estimate the bias as

$$\hat{B}_m = h^{p+1} \hat{m}^{(p+1)} \frac{1}{(p+1)!} e_0' \Gamma_p^{-1} \Lambda_{p,1}, \quad \text{where} \quad \hat{m}^{(p+1)} = [(p+1)!] e_{p+1}' H_q \Gamma_q^{-1} R_q' W_q Y / n.$$

The bias corrected estimator can then be written

$$\begin{aligned} \hat{m} - \hat{B}_m &= e_0' H_p \Gamma_p^{-1} R_p' W_p Y / n - h^{p+1} e_0' \Gamma_p^{-1} \Lambda_{p,1} e_{p+1}' H_q \Gamma_q^{-1} R_q' W_q Y / n \\ &= e_0' \Gamma_p^{-1} (R_p' W_p - \rho^{p+1} \Lambda_{p,1} e_{p+1}' \Gamma_q^{-1} R_q' W_q) Y / n, \end{aligned}$$

using the fact that  $e_{p+1}' H_q = \rho^{p+1} e_{p+1}'$ .

The fixed- $n$  variances are

$$\sigma_{\text{us}}^2 := (nh) \mathbb{V}[\hat{m}|X_1, \dots, X_n] = e_0' \Gamma_p^{-1} (h R_p' W_p \Sigma W_p R_p / n) \Gamma_p^{-1} e_0 \quad (20)$$

and

$$\begin{aligned} \sigma_{\text{rbc}}^2 &:= (nh) V[\hat{m} - \hat{B}_m | X_1, \dots, X_n] \\ &= e_0' \Gamma_p^{-1} (h/n) (R_p' W_p - \rho^{p+2} \Lambda_{p,1} e_{p+1}' \Gamma_q^{-1} R_q' W_q) \Sigma (R_p' W_p - \rho^{p+2} \Lambda_{p,1} e_{p+1}' \Gamma_q^{-1} R_q' W_q)' \Gamma_p^{-1} e_0, \end{aligned} \quad (21)$$

where

$$\Sigma = \text{diag}(v(X_i) : i = 1, \dots, n), \quad \text{with} \quad v(x) = \mathbb{V}[Y|X = x].$$

These are the closest analogue to the density case, but are still random due to the conditioning on the covariates. Their respective estimators are

$$\hat{\sigma}_{\text{us}}^2 = e_0' \Gamma_p^{-1} \left( h R_p' W_p \hat{\Sigma}_p W_p R_p \Gamma_p^{-1} / n \right) e_0$$

and

$$\hat{\sigma}_{\text{rbc}}^2 = e_0' \Gamma_p^{-1} (h/n) (R_p' W_p - \rho^{p+2} \Lambda_{p,1} e_{p+1}' \Gamma_q^{-1} R_q' W_q) \hat{\Sigma}_q (R_p' W_p - \rho^{p+2} \Lambda_{p,1} e_{p+1}' \Gamma_q^{-1} R_q' W_q)' \Gamma_p^{-1} e_0.$$

The conditional variance matrixes are estimated as

$$\hat{\Sigma}_p = \text{diag}(\hat{v}(X_i) : i = 1, \dots, n), \quad \text{with} \quad \hat{v}(X_i) = (Y_i - r_p(X_i - x)' \hat{\beta}_p)^2,$$

and

$$\hat{\Sigma}_q = \text{diag}(\hat{v}(X_i) : i = 1, \dots, n), \quad \text{with} \quad \hat{v}(X_i) = (Y_i - r_q(X_i - x)' \hat{\beta}_q)^2.$$

The Studentized statistics of interest are then:

$$T_{\text{us}} = \frac{\sqrt{nh}(\hat{m} - m)}{\hat{\sigma}_{\text{us}}}, \quad T_{\text{bc}} = \frac{\sqrt{nh}(\hat{m} - \hat{B}_m - m)}{\hat{\sigma}_{\text{us}}}, \quad T_{\text{rbc}} = \frac{\sqrt{nh}(\hat{m} - \hat{B}_m - m)}{\hat{\sigma}_{\text{rbc}}}.$$

The main result of this section is an Edgeworth expansion of the distribution function of these statistics.

### S.II.1.2 Edgeworth Expansion Terms

The terms of the Edgeworth expansion require further notation and discussion. The expressions are not nearly as compact as in the density case (cf. Section S.I.6).

Define the expectations of  $\Gamma_p$ ,  $\Gamma_q$ ,  $\Lambda_{p,k}$ , and  $\Lambda_{q,k}$  as  $\tilde{\Gamma}_p$ ,  $\tilde{\Gamma}_q$ ,  $\tilde{\Lambda}_{p,k}$ , and  $\tilde{\Lambda}_{q,k}$ , such as

$$\tilde{\Gamma}_p = \mathbb{E}[\Gamma_p] = \mathbb{E}[h^{-1}(K r_p r_p')(X_{h,i})].$$

These will be used to define nonrandom biases and variances that appear in the expansions.

The biases are defined in Eqn. (22), and are given by

$$\begin{aligned}\eta_{\text{us}} &= \sqrt{nh} \int e_0' \tilde{\Gamma}_p^{-1} K(u) r_p(u) (m(x - uh) - r_p(uh)' \beta_p) f(x - uh) du, \\ \eta_{\text{bc}} &= \sqrt{nh} \int e_0' \tilde{\Gamma}_p^{-1} K(u) r_p(u) (m(x - uh) - r_{p+1}(uh)' \beta_{p+1}) f(x - uh) du \\ &\quad - \sqrt{nh} \rho^{p+1} \int e_0' \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e_{p+1}' \tilde{\Gamma}_q^{-1} L(u) r_q(u) (m(x - ub) - r_q(ub)' \beta_q) f(x - ub) du.\end{aligned}$$

Further discussion and leading terms are found in Section S.II.4.

The fixed- $n$  variances are computed conditionally, and we must replace them with their nonrandom analogues (just as  $\eta_{\text{us}}$  and  $\eta_{\text{bc}}$  must be nonrandom). Recalling Equations (20) and (21), define

$$\tilde{\sigma}_{\text{us}}^2 := e_0' \tilde{\Gamma}_p^{-1} \tilde{\Psi}_p \tilde{\Gamma}_p^{-1} e_0,$$

where

$$\tilde{\Psi}_p = \mathbb{E} [\check{\Psi}_p] \quad \text{and} \quad \check{\Psi}_p := h R_p' W_p \Sigma W_p R_p / n,$$

and

$$\tilde{\sigma}_{\text{rbc}}^2 := e_0' \tilde{\Gamma}_p^{-1} \tilde{\Psi}_q \tilde{\Gamma}_p^{-1} e_0$$

where

$$\tilde{\Psi}_q = \mathbb{E} [\check{\Psi}_q] \quad \text{and} \quad \check{\Psi}_q := h \left( R_p' W_p - \rho^{p+2} \tilde{\Lambda}_{p,1} \tilde{\Gamma}_q^{-1} R_q' W_q \right) \Sigma \left( R_p' W_p / n - \rho^{p+2} \tilde{\Lambda}_{p,1} \tilde{\Gamma}_q^{-1} R_q' W_q / n \right)'.$$

In the course of the proofs, we will also use  $\hat{\Psi}_p = h R_p' W_p \hat{\Sigma}_p W_p R_p / n$  and the analogously-defined  $\hat{\Psi}_q$ .

We now give the precise forms of the polynomials in the Edgeworth expansion. As with the density, there will be both even and odd polynomials. These are not as compact or simple as the density case. Further, we will not attempt to simplify these functions by making use of limiting versions of moments. For example, we will *not* replace  $\tilde{\Lambda}_{p,1}$  by  $f(x) \int (K r_p)(u) u^{p+1} du$ , and similarly for other pieces. The only simplification made will be the use of  $q_{k,\text{us}}(z)$  in the expansion for  $T_{\text{bc}}$ , which otherwise would require further notation than what is below (along the lines of  $p_{1,\text{us}}(z)$  below).

First, define the following functions, which depend on  $n, p, q, h, b, K$  and  $L$ , but this is

generally suppressed:

$$\begin{aligned}
\ell_{\text{us}}^0(X_i) &= e'_0 \tilde{\Gamma}_p^{-1}(Kr_p)(X_{h,i}); \\
\ell_{\text{bc}}^0(X_i) &= \ell_{\text{us}}^0(X_i) - \rho^{p+2} e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1}(Lr_q)(X_{b,i}); \\
\ell_{\text{us}}^1(X_i, X_j) &= e'_0 \tilde{\Gamma}_p^{-1} \left( \mathbb{E}[(Kr_p r'_p)(X_{h,j})] - (Kr_p r'_p)(X_{h,j}) \right) \tilde{\Gamma}_p^{-1}(Kr_p)(X_{h,i}); \\
\ell_{\text{bc}}^1(X_i, X_j) &= \ell_{\text{us}}^1(X_i, X_j) - \rho^{p+2} e'_0 \tilde{\Gamma}_p^{-1} \left\{ \left( \mathbb{E}[(Kr_p r'_p)(X_{h,j})] - (Kr_p r'_p)(X_{h,j}) \right) \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e'_{p+1} \right. \\
&\quad \left. + ((Kr_p)(X_{h,j}) X_{h,i}^{p+1} - \mathbb{E}[(Kr_p)(X_{h,j}) X_{h,i}^{p+1}]) e'_{p+1} \right. \\
&\quad \left. + \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} \left( \mathbb{E}[(Lr_q r'_q)(X_{b,j})] - (Lr_q r'_q)(X_{b,j}) \right) \right\} \tilde{\Gamma}_q^{-1}(Lr_q)(X_{b,i}).
\end{aligned}$$

With this notation, we can write

$$\begin{aligned}
\tilde{\sigma}_{\text{us}}^2 &= \mathbb{E}[h^{-1} \ell_{\text{us}}^0(X)^2 v(X)], \\
\tilde{\sigma}_{\text{bc}}^2 &= \mathbb{E}[h^{-1} \ell_{\text{bc}}^0(X)^2 v(X)], \\
\eta_{\text{us}} &= s \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i) [m(X_i) - r_p(X_i - x)' \beta_p] \right],
\end{aligned}$$

and

$$\begin{aligned}
\eta_{\text{bc}} &= s \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i) [m(X_i) - r_{p+1}(X_i - x)' \beta_{p+1}] \right. \\
&\quad \left. - h^{-1} (\ell_{\text{bc}}^0(X_i) - \ell_{\text{us}}^0(X_i)) [m(X_i) - r_q(X_i - x)' \beta_q] \right].
\end{aligned}$$

We will define the Edgeworth expansion polynomials first for the undersmoothing case. The standard Normal density is  $\phi(z)$ . First, the even polynomials are

$$p_{1,\text{us}}(z) = \phi(z) \tilde{\sigma}_{\text{us}}^{-3} \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i)^3 \varepsilon_i^3 \right] \{ (2z^2 - 1)/6 \}$$

and

$$p_{3,\text{us}}(z) = -\phi(z) \tilde{\sigma}_{\text{us}}^{-1}.$$

The absence of  $p^{(2)}(z)$  is noteworthy: there is no version of this term for local polynomial estimation, because  $\varepsilon_i$  is conditionally mean zero.

Next, the odd polynomials for undersmoothing are defined as follows:

$$\begin{aligned}
q_{1,\text{us}}(z) &= \phi(z) \tilde{\sigma}_{\text{us}}^{-6} \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i)^3 \varepsilon_i^3 \right]^2 \{ z^3/3 + 7z/4 + \tilde{\sigma}_{\text{us}}^2 z(z^2 - 3)/4 \} \\
&\quad + \phi(z) \tilde{\sigma}_{\text{us}}^{-2} \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i) \ell_{\text{us}}^1(X_i, X_i) \varepsilon_i^2 \right] \{ -z(z^2 - 3)/2 \}
\end{aligned}$$



$$\begin{aligned}
& + \phi(z) \tilde{\sigma}_{\text{us}}^{-4} \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i)^4 (\varepsilon_i^4 - v(X_i)^2) \right] \{z(z^2 - 3)/8\} \\
& - \phi(z) \tilde{\sigma}_{\text{us}}^{-2} \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i)^2 r_p(X_{h,i})' \tilde{\Gamma}_p^{-1}(K r_p)(X_{h,i}) \varepsilon_i^2 \right] \{z(z^2 - 1)/2\} \\
& - \phi(z) \tilde{\sigma}_{\text{us}}^{-4} \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i)^3 r_p(X_{h,i})' \tilde{\Gamma}_p^{-1} \varepsilon_i^3 \right] \mathbb{E} \left[ h^{-1} (K r_p)(X_{h,i}) \ell_{\text{us}}^0(X_i) \varepsilon_i^2 \right] \{z(z^2 - 1)\} \\
& + \phi(z) \tilde{\sigma}_{\text{us}}^{-2} \mathbb{E} \left[ h^{-2} \ell_{\text{us}}^0(X_i)^2 (r_p(X_{h,i})' \tilde{\Gamma}_p^{-1}(K r_p)(X_{h,j}))^2 \varepsilon_j^2 \right] \{z(z^2 - 1)/4\} \\
& + \phi(z) \tilde{\sigma}_{\text{us}}^{-4} \mathbb{E} \left[ h^{-3} \ell_{\text{us}}^0(X_j)^2 r_p(X_{h,j})' \tilde{\Gamma}_p^{-1}(K r_p)(X_{h,i}) \ell_{\text{us}}^0(X_i) r_p(X_{h,j})' \tilde{\Gamma}_p^{-1}(K r_p)(X_{h,k}) \ell_{\text{us}}^0(X_k) \varepsilon_i^2 \varepsilon_k^2 \right] \\
& \quad \times \{z(z^2 - 1)/2\} \\
& + \phi(z) \tilde{\sigma}_{\text{us}}^{-4} \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i)^4 \varepsilon_i^4 \right] \{-z(z^2 - 3)/24\} \\
& + \phi(z) \tilde{\sigma}_{\text{us}}^{-4} \mathbb{E} \left[ h^{-1} (\ell_{\text{us}}^0(X_i)^2 v(X_i) - \mathbb{E}[\ell_{\text{us}}^0(X_i)^2 v(X_i)]) \ell_{\text{us}}^0(X_i)^2 \varepsilon_i^2 \right] \{z(z^2 - 1)/4\} \\
& + \phi(z) \tilde{\sigma}_{\text{us}}^{-4} \mathbb{E} \left[ h^{-2} \ell_{\text{us}}^1(X_i, X_j) \ell_{\text{us}}^0(X_i) \ell_{\text{us}}^0(X_j)^2 \varepsilon_j^2 v(X_i) \right] \{z(z^2 - 3)\} \\
& + \phi(z) \tilde{\sigma}_{\text{us}}^{-4} \mathbb{E} \left[ h^{-2} \ell_{\text{us}}^1(X_i, X_j) \ell_{\text{us}}^0(X_i) (\ell_{\text{us}}^0(X_j)^2 v(X_j) - \mathbb{E}[\ell_{\text{us}}^0(X_j)^2 v(X_j)]) \varepsilon_i^2 \right] \{-z\} \\
& + \phi(z) \tilde{\sigma}_{\text{us}}^{-4} \mathbb{E} \left[ h^{-1} (\ell_{\text{us}}^0(X_i)^2 v(X_i) - \mathbb{E}[\ell_{\text{us}}^0(X_i)^2 v(X_i)])^2 \right] \{-z(z^2 + 1)/8\};
\end{aligned}$$

$$q_{2,\text{us}}(z) = -\phi(z) \tilde{\sigma}_{\text{us}}^{-2} z/2;$$

$$q_{3,\text{us}}(z) = \phi(z) \tilde{\sigma}_{\text{us}}^{-4} \mathbb{E}[h^{-1} \ell_{\text{us}}^0(X_i)^3 \varepsilon_i^3] (z^3/3).$$

For robust bias correction, both the even polynomials,  $p_{1,\text{rbc}}(z)$  and  $p_{3,\text{rbc}}(z)$ , and the odd polynomials,  $q_{1,\text{rbc}}(z)$ ,  $q_{2,\text{rbc}}(z)$ , and  $q_{3,\text{rbc}}(z)$  are defined in the exact same way, but changing the  $\tilde{\sigma}_{\text{us}}$  to  $\tilde{\sigma}_{\text{rbc}}$ ,  $\ell_{\text{us}}^k(\cdot)$  to  $\ell_{\text{bc}}^k(\cdot)$ ,  $K$  to  $L$ , and  $p$  to  $q$ , and so forth.

The polynomials defined here are for *distribution function* expansions, and are different from those used for *coverage error*. The polynomials  $q_{1,\text{us}}$ ,  $q_{2,\text{us}}$ , and  $q_{3,\text{us}}$  and  $q_{1,\text{rbc}}$ ,  $q_{2,\text{rbc}}$ , and  $q_{3,\text{rbc}}$ , which do *not* have an argument, used for *coverage error* in the main text and in Corollary 8 below, are defined in terms of those given above, which *do* have an argument. Specifically, the polynomials above should be doubled, divided by the standard Normal density, and evaluated at the Normal quantile  $z_{\alpha/2}$ , that is,

$$q_{k,\bullet} := \frac{2}{\phi(z)} q_{k,\bullet}(z) \Big|_{z=z_{\alpha/2}}, \quad k = 1, 2, 3, \quad \bullet = \text{us}, \text{rbc}$$

For traditional bias correction,  $q_{1,\text{us}}(z)$ ,  $q_{2,\text{us}}(z)$ , and  $q_{3,\text{us}}(z)$  are used, but such simplification can not be done for  $p_{1,\text{bc}}(z)$  and  $p_{3,\text{bc}}(z)$ , which must be defined as

$$\begin{aligned}
p_{1,\text{bc}}(z) &= \phi(z) \tilde{\sigma}_{\text{us}}^{-3} \left( \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i)^3 \varepsilon_i^3 \right] \{-(z^2 - 1)/6\} + \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i)^2 \ell_{\text{bc}}^0(X_i) \varepsilon_i^3 \right] \{-(z^2 - 3)/4\} \right) \\
&+ \phi(z) \tilde{\sigma}_{\text{us}}^2 \tilde{\sigma}_{\text{rbc}}^{-5} \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_i)^2 \ell_{\text{bc}}^0(X_i) \varepsilon_i^3 \right] \{3(z^2 - 1)/4\}
\end{aligned}$$

and

$$p_{3,\text{bc}}(z) = -\phi(z)\tilde{\sigma}_{\text{us}}^{-1}.$$

Lastly, traditional bias correction also exhibits additional terms in the expansion (see discussion in the main text) representing the covariance of  $\hat{m}$  and  $\hat{B}_m$  (denoted by  $\Omega_{1,\text{bc}}$ ) and the variance of  $\hat{B}_m$  ( $\Omega_{2,\text{bc}}$ ). We now state their precise forms. These arise from the mismatch between the variance of the numerator of  $T_{\text{bc}}$  and the standardization used,  $\sigma_{\text{us}}^2$ , but these are random, and so  $\Omega_{1,\text{bc}}$  and  $\Omega_{2,\text{bc}}$  must be derived from the nonrandom versions,  $\tilde{\sigma}_{\text{rbc}}^2$  and  $\tilde{\sigma}_{\text{us}}^2$  (cf. Section S.I.6; for the same reason  $\eta_{\text{us}}$  and  $\eta_{\text{bc}}$  must be nonrandom). Recalling the definitions above,

$$\begin{aligned} \frac{\tilde{\sigma}_{\text{rbc}}^2}{\tilde{\sigma}_{\text{us}}^2} &= \frac{\mathbb{E}[h^{-1}\ell_{\text{bc}}^0(X)^2v(X)]}{\mathbb{E}[h^{-1}\ell_{\text{us}}^0(X)^2v(X)]} \\ &= \frac{\mathbb{E}[h^{-1}\{\ell_{\text{us}}^0(X) + (\ell_{\text{bc}}^0(X) - \ell_{\text{us}}^0(X))\}^2v(X)]}{\mathbb{E}[h^{-1}\ell_{\text{us}}^0(X)^2v(X)]} \\ &= 1 - 2\tilde{\sigma}_{\text{us}}^{-2}\mathbb{E}[h^{-1}\{\ell_{\text{us}}^0(X)(\ell_{\text{bc}}^0(X) - \ell_{\text{us}}^0(X))\}v(X)] + \tilde{\sigma}_{\text{us}}^{-2}\mathbb{E}[h^{-1}\{(\ell_{\text{bc}}^0(X) - \ell_{\text{us}}^0(X))^2v(X)] \\ &= 1 - 2\rho^{1+(p+1)}\tilde{\sigma}_{\text{us}}^{-2}\mathbb{E}[h^{-1}\{\rho^{-p-2}\ell_{\text{us}}^0(X)(\ell_{\text{bc}}^0(X) - \ell_{\text{us}}^0(X))\}v(X)] \\ &\quad + \rho^{1+2(p+1)}\tilde{\sigma}_{\text{us}}^{-2}\mathbb{E}[h^{-1}\{\rho^{-p-2}(\ell_{\text{bc}}^0(X) - \ell_{\text{us}}^0(X))^2v(X)] \end{aligned}$$

Therefore

$$\Omega_{1,\text{bc}} = -2\tilde{\sigma}_{\text{us}}^{-2}\mathbb{E}[h^{-1}\{\rho^{-p-2}\ell_{\text{us}}^0(X)(\ell_{\text{bc}}^0(X) - \ell_{\text{us}}^0(X))\}v(X)]$$

and

$$\Omega_{2,\text{bc}} = \tilde{\sigma}_{\text{us}}^{-2}\mathbb{E}[h^{-1}\{\rho^{-p-2}(\ell_{\text{bc}}^0(X) - \ell_{\text{us}}^0(X))^2v(X)].$$

**Remark 7** (Simplifications). It is possible for the above-defined polynomials to simplify in special cases. A leading example is in the homoskedastic Gaussian regression model:

$$Y_i = m(X_i) + \varepsilon_i, \quad \text{where} \quad \varepsilon_i \sim \mathcal{N}(0, v).$$

This model is a common theoretical baseline to study, though over-simplified from an empirical point of view. In this special case,  $\mathbb{E}[\varepsilon_i^3] = 0$  and thus  $q_{3,\text{us}}(z) \equiv 0$ , entirely removing this term from the Edgeworth expansions. This has little bearing on the conceptual conclusions however, and in particular the comparison of undersmoothing and robust bias correction.

## S.II.2 Details of Practical Implementation

In the main text we give a direct plug-in (DPI) rule to implement the coverage-error optimal bandwidth. Here we give complete details for this procedure as well as document a second practical choice, based on a rule-of-thumb (ROT) strategy. Both choices yield the optimal coverage error decay rate at interior and boundary points. All our methods are implemented in software available from the authors' websites and via the R package `nprobust` available at <https://cran.r-project.org/package=nprobust>.

As in the density case, the MSE-optimal bandwidth undercovers when used in the under-smoothing confidence interval; that is, Remark 1 applies directly. See also Hall and Horowitz (2013).

### S.II.2.1 Bandwidth Choice: Rule-of-Thumb (ROT)

As with the density case, a simple rule-of-thumb based on rescaling the MSE-optimal bandwidth is:

$$\hat{h}_{\text{rot}}^{\text{int}} = \hat{h}_{\text{mse}}^{\text{int}} n^{-(p-1)/((2p+3)(p+4))} \quad \text{and} \quad \hat{h}_{\text{rot}}^{\text{bnd}} = \hat{h}_{\text{mse}}^{\text{bnd}} n^{-p/((2p+3)(p+3))}.$$

where  $\hat{h}_{\text{mse}}^{\text{int}}$  and  $\hat{h}_{\text{mse}}^{\text{bnd}}$  denote readily-available implementations of the MSE-optimal bandwidth for interior and boundary points, respectively. See, e.g., Fan and Gijbels (1996). Again, when  $p = 1$  in the interior, no scaling is needed ( $\hat{h}_{\text{rot}}^{\text{int}} = \hat{h}_{\text{mse}}^{\text{int}}$ ), but for  $p > 1$  any data-driven MSE-optimal bandwidth should always be shrunk to improve inference at the boundary (i.e., reduce coverage errors of the robust bias-corrected confidence intervals).

The ROT selector may be especially attractive for simplicity, if estimating the constants described below in the DPI case is prohibitive.

Remark 2 applies to this case as well, though less transparently and without consequences that are as dramatic.

### S.II.2.2 Bandwidth Choice: Direct Plug-In (DPI)

We now detail the required steps to implement the plug-in bandwidth  $\hat{h}_{\text{dpi}}$  for interior and boundary points. We always set  $K = L$ ,  $\rho = 1$ , and  $q = p + 1$ . The steps are:

- (1) As a pilot bandwidth, use  $\hat{h}_{\text{mse}}$ : any data-driven version of  $h_{\text{mse}}^*$ .

- (2) Using this bandwidth, estimate the regression function  $m(X_i)$  as  $\hat{m}(X_i; \hat{h}_{\text{mse}}) = r_p(X_i - x)' \hat{\beta}_p(\hat{h}_{\text{mse}})$ , where  $\hat{\beta}_p(\hat{h}_{\text{mse}})$  is the local polynomial coefficient estimate of order  $p$  exactly as defined in the main text, using the bandwidth  $\hat{h}_{\text{mse}}$ .

Form  $\hat{\varepsilon}_i = Y_i - \hat{m}(X_i; \hat{h}_{\text{mse}})$ .

- (3) Following [Fan and Gijbels \(1996, §4.2\)](#) we estimate derivatives  $m^{(k)}$  using a global least squares polynomial fit of order  $k + 2$ . That is, estimate  $\hat{m}^{(p+3)}(x)$  as

$$\hat{m}^{(p+3)}(x) = [\hat{\gamma}]_{p+4} (p+3)! + [\hat{\gamma}]_{p+5} (p+4)! x + [\hat{\gamma}]_{p+6} \frac{(p+5)!}{2} x^2,$$

where  $[\hat{\gamma}]_k$  is the  $k$ -th element of the vector  $\hat{\gamma}$  that is estimated as

$$\hat{\gamma} = \arg \min_{\gamma \in \mathbb{R}^{p+6}} \sum_{i=1}^n (Y_i - r_{p+5}(X_i)' \gamma)^2.$$

The estimate for  $\hat{m}^{(p+2)}(x)$  is similar, with all indexes incremented down once.

For interior points, both are needed, while only  $\hat{m}^{(p+2)}(x)$  is required for the boundary.

- (4) The estimated polynomials  $\hat{q}_{k,\text{rbc}}$ ,  $k = 1, 2, 3$  and the bias constants  $\hat{\eta}_{\text{bc}}^{\text{int}}$  and  $\hat{\eta}_{\text{bc}}^{\text{bnd}}$  are defined as follows. The polynomials  $q_{1,\text{rbc}}$ ,  $q_{2,\text{rbc}}$ , and  $q_{3,\text{rbc}}$ , which do *not* have an argument, are defined in terms of those given in [Section S.II.1.2](#), which *do* have an argument. Specifically, the polynomials in [Section S.II.1.2](#) should be doubled, divided by the standard Normal density, and evaluated at the Normal quantile  $z_{\alpha/2}$ , that is,  $q_{k,\text{rbc}} = \phi(z_{\alpha/2})^{-1} q_{k,\text{rbc}}(z_{\alpha/2})$ .

Note that with the recommended choice of  $K = L$ ,  $\rho = 1$ , and  $q = p + 1$ , the polynomials  $\hat{q}_{k,\text{rbc}}$ ,  $k = 1, 2, 3$  can be read off the expressions for the undersmoothing versions,  $\hat{q}_{k,\text{us}}$ ,  $k = 1, 2, 3$ , with  $p$  replaced by  $p + 1$ .

The bias terms, for the interior and boundary, are given as follows. With  $q = p + 1$ , and hence even, and  $\rho = 1$ , the expressions of [Section S.II.4](#) simplify. For the interior:  $\eta_{\text{bc}}^{\text{int}} = \sqrt{nh} h^{p+3} \tilde{\eta}_{\text{bc}}^{\text{int}}$ , with

$$\begin{aligned} \tilde{\eta}_{\text{bc}}^{\text{int}} = & \frac{m^{(p+2)}}{(p+2)!} \left\{ e'_0 \tilde{\Gamma}_p^{-1} \left( \tilde{\Lambda}_{p,2} - \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} \tilde{\Lambda}_{q,1} \right) \right\} \\ & + \frac{m^{(p+3)}}{(p+3)!} \left\{ e'_0 \tilde{\Gamma}_p^{-1} \left( \tilde{\Lambda}_{p,3} - \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} \tilde{\Lambda}_{q,2} \right) \right\}; \end{aligned}$$

At the boundary:  $\eta_{bc}^{\text{bnd}} = \sqrt{nh}h^{p+2}\tilde{\eta}_{bc}^{\text{bnd}}$ , with

$$\tilde{\eta}_{bc}^{\text{bnd}} = \frac{m^{(p+2)}}{(p+2)!} \left\{ e'_0 \tilde{\Gamma}_p^{-1} \left( \tilde{\Lambda}_{p,2} - \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} \tilde{\Lambda}_{q,1} \right) \right\}.$$

The estimates of these,  $\hat{q}_{k,\text{rbc}}$ ,  $k = 1, 2, 3$  and  $\hat{\eta}_{bc}^{\text{int}}$  and  $\hat{\eta}_{bc}^{\text{bnd}}$ , are defined by replacing:

- (i)  $h$  with  $\hat{h}_{\text{mse}}$ ,
  - (ii) population expectations with sample averages (see note below),
  - (iii) residuals  $\varepsilon_i$  with  $\hat{\varepsilon}_i$ ,
  - (iv) derivatives  $m^{(p+2)}$  and  $m^{(p+3)}$  with their estimators from above,
  - (v) limiting matrixes  $\tilde{\Gamma}_p$ ,  $\tilde{\Lambda}_{p,2}$ , etc, with the corresponding sample versions using the bandwidth  $\hat{h}_{\text{mse}}$ , e.g.,  $\tilde{\Gamma}_p$  is replaced with  $\Gamma_p(\hat{h}_{\text{mse}}) = R'_p W_p(\hat{h}_{\text{mse}}) R_p / n$ , where  $W_p(\hat{h}_{\text{mse}}) = \text{diag} \left( \hat{h}_{\text{mse}}^{-1} K \left( (X_i - x) / \hat{h}_{\text{mse}} \right) \right)$ .
- (5) Finally  $\hat{h}_{\text{dpi}}^{\text{int}} = \hat{H}_{\text{dpi}}^{\text{int}}(\hat{h}_{\text{mse}})n^{-1/(p+4)}$  and  $\hat{h}_{\text{dpi}}^{\text{bnd}} = \hat{H}_{\text{dpi}}^{\text{bnd}}(\hat{h}_{\text{mse}})n^{-1/(p+3)}$ , where

$$\hat{H}_{\text{dpi}}^{\text{int}}(\hat{h}_{\text{mse}}) = \arg \min_H \left| H^{-1} \hat{q}_{1,\text{rbc}} + H^{1+2(p+3)} (\hat{\eta}_{bc}^{\text{int}})^2 \hat{q}_{2,\text{rbc}} + H^{p+3} (\hat{\eta}_{bc}^{\text{int}}) \hat{q}_{3,\text{rbc}} \right|,$$

while at (or near) the boundary the optimal bandwidth is  $h_{\text{rbc}}^* = H_{\text{rbc}}^*(\rho)n^{-1/(p+3)}$ , where

$$\hat{H}_{\text{dpi}}^{\text{bnd}}(\hat{h}_{\text{mse}}) = \arg \min_H \left| H^{-1} \hat{q}_{1,\text{rbc}} + H^{1+2(p+2)} (\hat{\eta}_{bc}^{\text{bnd}})^2 \hat{q}_{2,\text{rbc}} + H^{p+2} (\hat{\eta}_{bc}^{\text{bnd}}) \hat{q}_{3,\text{rbc}} \right|.$$

These numerical minimizations are easily solved; see note below. Code available from the authors' websites performs all the above steps.

**Remark 8** (Notes on computation).

- When numerically solving the above minimization problems, computation will be greatly sped up by squaring the objective function.
- For step 4(ii) above, in estimating  $\hat{q}_{1,\text{rbc}}$ , and specifically when replacing population expectations with sample averages, we use the appropriate  $U$ -statistic forms to reduce bias. There are several terms which are expectations over two or three observations, and for these the second or third order  $U$ -statistic forms are preferred. For example, when estimating terms such as

$$\mathbb{E} \left[ h^{-2} \ell_{\text{us}}^0(X_i)^2 (r_p(X_{h,i})' \tilde{\Gamma}_p^{-1} (K r_p)(X_{h,j}))^2 \varepsilon_j^2 \right]$$

we use

$$\frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i} \left[ \hat{h}_{\text{mse}}^{-2} \hat{\ell}_{\text{rbc}}^0(X_i)^2 (r_p(X_{\hat{h}_{\text{mse},i}})' \Gamma_p^{-1} (K r_p)(X_{\hat{h}_{\text{mse},j}}))^2 \hat{\varepsilon}_j^2 \right],$$

where  $\hat{\ell}_{\text{rbc}}^0(X_i)$  is made feasible as in step 4(v). ■

### S.II.2.3 Alternative Standard Errors

As argued in the main text, using variance forms other than (20) and (21) can be detrimental to coverage. Within these forms however, two alternative estimates of  $\Sigma$  are natural. First, motivated by the fact that the least-squares residuals are on average too small, the well-known HCK class of heteroskedasticity consistent estimators can be used; see [MacKinnon \(2013\)](#) for details and a recent review. In our notation, these are defined as follows. First,  $\hat{\sigma}_{\text{us}}^2$ -HCK0 is the estimator above. Then, for  $k = 1, 2, 3$ , the  $\hat{\sigma}_{\text{us}}^2$ -HCK estimator is obtained by dividing  $\hat{\varepsilon}_i^2$  by, respectively,  $(n - 2 \text{tr}(Q_p) + \text{tr}(Q_p' Q_p))/n$ ,  $(1 - Q_{p,ii})$ , and  $(1 - Q_{p,ii})^2$ , where  $Q_{p,ii}$  is the  $i$ -th diagonal element of the projection matrix  $Q_p := R_p' \Gamma_p^{-1} R_p' W_p / n$ . The corresponding estimators  $\hat{\sigma}_{\text{rbc}}^2$ -HCK are the same way, with  $q$  in place of  $p$ . As is well-known in the literature, these estimators perform better for small sample sizes, a fact we confirm in our simulation study below.

A second option is to use a nearest-neighbor-based variance estimators with a fixed number of neighbors, following the ideas of [Muller and Stadtmuller \(1987\)](#); [Abadie and Imbens \(2008\)](#). To define these, let  $J$  be a fixed number and  $j(i)$  be the  $j$ -th closest observation to  $X_i$ ,  $j = 1, \dots, J$ , and set  $\hat{v}(X_i) = \frac{J}{J+1} (Y_i - \sum_{j=1}^J Y_{j(i)} / J)^2$ . This “estimate” is unbiased (but inconsistent) for  $v(X_i)$ .

Both types of residual estimators could be handled in our results. The constants will change, but the rates will not. This is because, in all cases, the errors in estimating  $v(X_i)$  are no greater than in the original  $\hat{m}(x)$ . Inspection of the proof shows that simple modifications allow for the HCK estimators: only the terms of Eqn. (27) will change, and indeed, we conjecture that the HCK estimators will result in fewer terms and a reduced coverage error. This is consistent with the improved finite-sample behavior of these estimators and the fact that they are asymptotically equivalent. Accommodating the nearest-neighbor estimates require slightly more work and a modified version of Assumption S.II.3.3.

One crucial property of our method, in the context of Edgeworth expansions, is that the bias in estimation of  $\Sigma$  is of the same order as the original  $\hat{m}(x)$ . Using other methods may result in additional terms, with possibly distinct rates, appearing in the Edgeworth expansions.

Some examples that may have this issue are (i) using  $\hat{v}(X_i) = (Y_i - \hat{m}(x))^2$ ; (ii) using local or assuming global heteroskedasticity; (iii) using other nonparametric estimators for  $v(X_i)$ , relying on new tuning parameters.

## S.II.3 Assumptions

Copied directly from the main text (see discussion there), the following assumptions are sufficient for our results.

**Assumption S.II.3.1** (Data-generating process).  *$\{(Y_1, X_1), \dots, (Y_n, X_n)\}$  is a random sample, where  $X_i$  has the absolutely continuous distribution with Lebesgue density  $f$ ,  $\mathbb{E}[Y^{8+\delta}|X] < \infty$  for some  $\delta > 0$ , and in a neighborhood of  $x$ ,  $f$  and  $v$  are continuous and bounded away from zero,  $m$  is  $S > q + 2$  times continuously differentiable with bounded derivatives, and  $m^{(S)}$  is Hölder continuous with exponent  $\varsigma$ .*

**Assumption S.II.3.2** (Kernels). *The kernels  $K$  and  $L$  are positive, bounded, even functions, and with compact support.*

**Assumption S.II.3.3** (Cramér's Condition). *For each  $\delta > 0$  and all sufficiently small  $h$ , the random variables  $Z_{\text{us}}(u)$  and  $Z_{\text{rbc}}(u)$  defined below obey*

$$\sup_{t \in \mathbb{R}^{\dim\{Z(u)\}}, \|t\| > \delta} \left| \int \exp\{it'Z(u)\} f(x - uh) du \right| \leq 1 - C(x, \delta)h,$$

where  $C(x, \delta) > 0$  is a fixed constant,  $\|t\|^2 = \sum_{d=1}^{\dim\{Z(u)\}} t_d^2$ , and  $i = \sqrt{-1}$ .

The random variables of Assumption S.II.3.3 are defined follows. For two kernels  $K_1$  and  $K_2$ , two polynomial orders (i.e. positive integers)  $p_1$  and  $p_2$ , a bandwidth  $b$ , and a scalar  $\rho$ , let

$$Z_m(u; K_1, p_1, p_2, b, \rho) := \left( K_1(u)r_{p_1}(u)' \varepsilon, K_1(u)r_{p_1}(u)'(m(x-ub) - r_{p_2}(ub)' \beta_{p_2}), \text{vech}(K_1(u)r_{p_1}(u)r_{p_1}(u)') \right)$$

and

$$\begin{aligned} Z_\sigma(u; K_1, K_2, p_1, p_2, b, \rho) := & \left( \text{vech}(K_1(u)K_2(u\rho)r_{p_1}(u)r_{p_2}(u\rho)'\varepsilon^2)', \right. \\ & \text{vech}(K_1(u)K_2(u\rho)r_{p_1}(u)r_{p_2}(u\rho)'v(x-ub))', \\ & \text{vech}(K_1(u)K_2(u\rho)r_{p_1}(u)r_{p_2}(u\rho)'\varepsilon(m(x-ub) - r_{p_2}(ub)'\beta_{p_2}))', \\ & \left. \text{vech}(K_2(u)^2r_{p_2}(u)r_{p_2}(u)')' \right), \end{aligned}$$

$$\begin{aligned} & \text{vech}(K_1(u)K_2(u\rho)r_{p_1}(u)r_{p_2}(u\rho)'r_{p_2}(u)'\varepsilon)', \\ & \text{vech}(K_1(u)K_2(u\rho)r_{p_1}(u)r_{p_2}(u\rho)'r_{p_2}(u\rho)'\varepsilon(m(x-ub) - r_{p_2}(ub)'\beta_{p_2}))')' \end{aligned}$$

The subscripts are intended to make clear that  $Z_m(\cdot)$  collects quantities from the numerator of the Studentized statistic, while  $Z_\sigma(\cdot)$  gathers additional variables required for the variance estimation. With this notation, we define

$$Z_{\text{us}}(u) = (Z_m(u; K, p, p, h, 1)', Z_\sigma(u; K, K, p, p, h, 1)')',$$

$$Z_{\text{bc}}(u) = (Z_m(u; K, p, p+1, h, 1)', Z_m(u; L, q, q, b, \rho)', \text{vech}(K(u)r_p(u)u^{p+1})', Z_\sigma(u; K, K, p, p, h, 1)')',$$

and

$$\begin{aligned} Z_{\text{rbc}}(u) = & (Z_m(u; K, p, p+1, h, 1)', Z_m(u; L, q, q, b, \rho)', \text{vech}(K(u)r_p(u)u^{p+1})', \\ & Z_\sigma(u; K, K, p, q, b, \rho)', Z_\sigma(u; L, L, q, q, b, 1)', Z_\sigma(u; K, L, p, q, b, \rho)')'. \end{aligned}$$

This notation is quite compact, and while it emphasizes the simplicity of Cramér's condition and the fact that it puts mild restrictions on the kernels, it does obscure the full notational breadth, particularly for  $Z_{\text{rbc}}$ . This is mostly repetitive: what holds for the kernel  $K$  and order  $p$  fit must also hold for  $L$  and  $q$ , and for their squares and cross products. To make this clear, we can expand all the  $Z_m$  and  $Z_\sigma$ , to write out the full statistics

$$\begin{aligned} Z_{\text{us}}(u) = & \left( K(u)r_p(u)'\varepsilon, K(u)r_p(u)'(m(x-uh) - r_p(uh)'\beta_p), \text{vech}(K(u)r_p(u)r_p(u)')', \right. \\ & \text{vech}(K(u)^2r_p(u)r_p(u)'\varepsilon^2)', \text{vech}(K(u)^2r_p(u)r_p(u)'v(x-uh))', \\ & \text{vech}(K(u)^2r_p(u)r_p(u)'\varepsilon(m(x-uh) - r_p(uh)'\beta_p))', \text{vech}(K(u)^2r_p(u)r_p(u)'r_p(u)')', \\ & \left. \text{vech}(K(u)^2r_p(u)r_p(u)'r_p(u)'\varepsilon)', \text{vech}(K(u)^2r_p(u)r_p(u)'r_p(u)'\varepsilon(m(x-uh) - r_p(uh)'\beta_p))' \right)', \\ Z_{\text{bc}}(u) = & \left( K(u)r_p(u)'\varepsilon, \text{vech}(K(u)r_p(u)r_p(u)')', \right. \\ & \text{vech}(K(u)^2r_p(u)r_p(u)'\varepsilon^2)', \text{vech}(K(u)^2r_p(u)r_p(u)'v(x-uh))', \\ & \text{vech}(K(u)^2r_p(u)r_p(u)'\varepsilon(m(x-uh) - r_p(uh)'\beta_p))', \text{vech}(K(u)^2r_p(u)r_p(u)'r_p(u)')', \\ & \text{vech}(K(u)^2r_p(u)r_p(u)'r_p(u)'\varepsilon)', \text{vech}(K(u)^2r_p(u)r_p(u)'r_p(u)'\varepsilon(m(x-uh) - r_p(uh)'\beta_p))', \\ & K(u)r_p(u)'(m(x-uh) - r_{p+1}(uh)'\beta_{p+1}), L(u\rho)r_q(u\rho)'\varepsilon, \text{vech}(L(u\rho)r_q(u\rho)r_q(u\rho)')', \\ & \left. \text{vech}(K(u)r_p(u)u^{p+1})', L(u\rho)r_q(u\rho)'(m(x-uh) - r_q(uh)'\beta_q) \right)', \end{aligned}$$



and

$$\begin{aligned}
Z_{\text{rbc}}(u) = & \left( Z_{\text{bc}}(u)', \text{vech}(K(u)^2 r_p(u) r_p(u)' \varepsilon^2)', \text{vech}(K(u)^2 r_p(u) r_p(u)' v(x - ub))', \right. \\
& \text{vech}(K(u)^2 r_p(u) r_p(u)' \varepsilon(m(x - ub) - r_q(ub)' \beta_q))', \text{vech}(K(u)^2 r_p(u) r_p(u)' r_q(u \rho)')', \\
& \text{vech}(K(u)^2 r_p(u) r_p(u)' r_q(u \rho)' \varepsilon)', \text{vech}(K(u)^2 r_p(u) r_p(u)' r_q(u \rho)' \varepsilon(m(x - ub) - r_q(ub)' \beta_q))', \\
& \text{vech}(L(u)^2 r_q(u) r_q(u)' \varepsilon^2)', \text{vech}(L(u)^2 r_q(u) r_q(u)' v(x - ub))', \\
& \text{vech}(L(u)^2 r_q(u) r_q(u)' \varepsilon(m(x - ub) - r_q(ub)' \beta_q))', \text{vech}(L(u)^2 r_q(u) r_q(u)' r_q(u)')', \\
& \text{vech}(L(u)^2 r_q(u) r_q(u)' r_q(u)' \varepsilon)', \text{vech}(L(u)^2 r_q(u) r_q(u)' r_q(u)' \varepsilon(m(x - ub) - r_q(ub)' \beta_q))', \\
& \text{vech}(K(u) L(u \rho) r_p(u) r_q(u \rho)' \varepsilon^2)', \text{vech}(K(u) L(u \rho) r_p(u) r_q(u \rho)' v(x - ub))', \\
& \text{vech}(K(u) L(u \rho) r_p(u) r_q(u \rho)' \varepsilon(m(x - ub) - r_q(ub)' \beta_q))', \text{vech}(L(u)^2 r_q(u) r_q(u)' r_q(u)')', \\
& \text{vech}(K(u) L(u \rho) r_p(u) r_q(u \rho)' r_q(u)' \varepsilon)', \\
& \left. \text{vech}(K(u) L(u \rho) r_p(u) r_q(u \rho)' r_q(u \rho)' \varepsilon(m(x - ub) - r_q(ub)' \beta_q))' \right)'.
\end{aligned}$$

**Remark 9** (Sufficient Conditions for Cramér's Condition). Assumption [S.II.3.3](#) is a high level condition, but one that is fairly mild. It is essentially a continuity requirement, and is discussed at length by (among others) [Bhattacharya and Rao \(1976\)](#), [Bhattacharya and Ghosh \(1978\)](#), and [Hall \(1992a\)](#). For a recent work in econometrics, the present condition can be compared to that employed by [Kline and Santos \(2012\)](#) for parametric regression (the role of the covariates is here played by  $r_p(X_{h,i})$ ): ours is more complex due to the nonparametric smoothing bias and the fact that the expansion is carried out to higher order.

It is straightforward to provide sufficient conditions for Assumption [S.II.3.3](#), given that Assumptions [S.II.3.1](#) and [S.II.3.2](#) hold. In particular, if we additionally assume that

$$(1, \text{vech}(K(u) r_p(u) r_p(u)'), \text{vech}(K(u)^2 r_p(u) r_p(u)' r_p(u)'))'$$

comprises a linearly independent set of functions on  $[-1, 1]$ , then it holds  $Z_{\text{us}}(u)$  has components that are nondegenerate and absolutely continuous, and this will imply that Assumption [S.II.3.3](#) holds for  $Z_{\text{us}}(u)$ , by arguing as in [Bhattacharya and Ghosh \(1978, Lemma 2.2\)](#) and [Hall \(1992a, p. 65\)](#). This is precisely the approach taken by [Chen and Qin \(2002\)](#), when studying undersmoothed local linear regression. If the linear independence continues to hold when the set of functions is augmented with  $\text{vech}(L(u) r_q(u) r_q(u)')$ , then  $Z_{\text{bc}}(u)$  satisfies Assumption [S.II.3.3](#) as well. To obtain the result for  $Z_{\text{rbc}}(u)$  requires that linear independence hold for

$$(1, \text{vech}(K(u) r_p(u) r_p(u)'), \text{vech}(K(u)^2 r_p(u) r_p(u)' r_q(u)')', \text{vech}(L(u) r_q(u) r_q(u)'),$$

$$\text{vech}(L(u)^2 r_q(u) r_q(u)' r_q(u)')', \text{vech}(K(u) L(u \rho) r_p(u) r_q(u \rho)' r_q(u \rho)')').$$

At heart, these are requirements on the kernel functions, just as in Assumption [S.I.3.3](#) in the density case. The uniform kernel is again ruled out. See Remark [5](#). Further, note that if these sets of functions are not linearly independent, there will exist a smaller set of functions which are linearly independent and can replace the original set while leaving the value of the statistic unchanged (see [Bhattacharya and Ghosh \(1978, p. 442\)](#)).

In sum, this makes clear that Assumption [S.II.3.3](#) is quite mild. ■

Finally, the precise random variables  $Z_{\text{us}}(u)$ ,  $Z_{\text{bc}}(u)$ , and  $Z_{\text{rbc}}(u)$  used can be replaced with slightly different constructions without altering the conclusions of Theorem [5](#): there are other potential functions  $\tilde{T}$  that satisfy Eqn. [\(23\)](#) in the proof. Such changes necessarily involve asymptotically negligible terms, and do not materially alter the severity of the restrictions imposed.

## S.II.4 Bias

We will not present a detailed discussion of bias issues, along the lines of Section [S.I.4.1](#), for brevity; we focus only on the case of nonbinding smoothness.

The biases  $\eta_{\text{us}}$  and  $\eta_{\text{bc}}$  are not as conceptually simple as in the density case. The closest parallel to the density case would be (for example)  $\eta_{\text{us}} = \sqrt{nh}(\mathbb{E}[\hat{m}] - m)$ , but this can not be used due to the presence of  $\Gamma_p^{-1}$  inside the expectation, and the next natural choice, the conditional bias  $\sqrt{nh}(\mathbb{E}[\hat{m}|X_1, \dots, X_n] - m)$ , is still random. Instead,  $\eta_{\text{us}}$  and  $\eta_{\text{bc}}$  are biases computed after replacing  $\Gamma_p$ ,  $\Gamma_q$ , and  $\Lambda_{p,1}$  with their expectations, denoted  $\tilde{\Gamma}_p$ ,  $\tilde{\Gamma}_q$ , and  $\tilde{\Lambda}_{p,1}$ . We thus define

$$\begin{aligned} \eta_{\text{us}} &= \sqrt{nh} \int e_0' \tilde{\Gamma}_p^{-1} K(u) r_p(u) (m(x - uh) - r_p(uh)' \beta_p) f(x - uh) du, \\ \eta_{\text{bc}} &= \sqrt{nh} \int e_0' \tilde{\Gamma}_p^{-1} K(u) r_p(u) (m(x - uh) - r_{p+1}(uh)' \beta_{p+1}) f(x - uh) du \\ &\quad - \sqrt{nh} \rho^{p+1} \int e_0' \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e_{p+1}' \tilde{\Gamma}_q^{-1} L(u) r_q(u) (m(x - ub) - r_q(ub)' \beta_q) f(x - ub) du. \end{aligned} \tag{22}$$

For the generic results of coverage error or the generic Edgeworth expansions of Theorem [5](#) below, the above definitions of  $\eta_{\text{us}}$  and  $\eta_{\text{bc}}$  are suitable. For the Corollaries detailing specific cases, and to understand the behavior at different points, it is useful to make the leading terms precise, that is, analogues of Equations [\(10\)](#) and [\(11\)](#). We must consider interior and boundary

point estimation, and even and odd  $q$ . We depart slightly from other terms of the expansion in that we do retain only the leading term for some pieces. This is done in order to capture the rate of convergence explicitly and to give practicable results. These results are derived by [Fan and Gijbels \(1996, Section 3.7\)](#) and similar calculations (though our expressions differ slightly as fixed- $n$  expectations are retained as much as possible).

Since  $p$  is odd, both at boundary and interior points we have

$$\eta_{\text{us}} = \sqrt{nh} h^{p+1} \frac{m^{(p+1)}}{(p+1)!} e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} [1 + o(1)].$$

Moving to  $\eta_{\text{bc}}$ , consider the first term, which in the present notation is:  $\sqrt{nh} \mathbb{E}[h^{-1} \ell_{\text{bc}}^0(X)(m(X) - r_{p+1}(X - x)' \beta_{p+1})]$ . With  $p+1$  even, we find that in the interior the leading terms are

$$\sqrt{nh} h^{p+3} e'_0 \tilde{\Gamma}_p^{-1} \left( \frac{m^{(p+2)}}{(p+2)!} \tilde{\Lambda}_{p,2} + \frac{m^{(p+3)}}{(p+3)!} \tilde{\Lambda}_{p,3} \right) [1 + o(1)],$$

due to the well-known symmetry properties of local polynomials that result in the cancellation of the leading terms of  $\tilde{\Gamma}_p^{-1}$  and  $\tilde{\Lambda}_{p,2}$ . The rate of  $h^{p+3}$  accounts for this. At the boundary, no such cancellation occurs and we have only

$$\sqrt{nh} h^{p+2} \frac{m^{(p+2)}}{(p+2)!} e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,2} [1 + o(1)].$$

Next, turn to the bias of the bias estimate:

$$\sqrt{nh} \rho^{p+1} e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} \int L(u) r_q(u) (m(x - ub) - r_q(ub)' \beta_q) f(x - ub) du.$$

If  $q$  is odd (so that  $q - (p+1)$  is also odd), then at the interior or boundary the leading term will be

$$\sqrt{nh} b^{q+1} \rho^{p+1} \frac{m^{(q+1)}}{(q+1)!} e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} \tilde{\Lambda}_{q,1} [1 + o(1)] \asymp \sqrt{nh} h^{p+1} b^{q-p}.$$

The same expression applies at the boundary for  $q$  even. However, for the interior, if  $q$  is even, then we again have cancellation of certain leading terms, resulting in the bias of the bias estimate being

$$\sqrt{nh} b^{q+2} \rho^{p+1} e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} \left( \frac{m^{(q+1)}}{(q+1)!} \tilde{\Lambda}_{q,1} + \frac{m^{(q+2)}}{(q+2)!} \tilde{\Lambda}_{q,2} \right) [1 + o(1)] \asymp \sqrt{nh} h^{p+1} b^{q+1-p}.$$

Combining all these results, we find the following (dropping remainder terms): for an interior

point, with  $q$  even,

$$\eta_{\text{bc}} = \sqrt{nh}h^{p+3} \left\{ e'_0 \tilde{\Gamma}_p^{-1} \left( \frac{m^{(p+2)}}{(p+2)!} \tilde{\Lambda}_{p,2} + \frac{m^{(p+3)}}{(p+3)!} \tilde{\Lambda}_{p,3} \right) - \rho^{-2} b^{q-(p+1)} e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} \left( \frac{m^{(q+1)}}{(q+1)!} \tilde{\Lambda}_{q,1} + \frac{m^{(q+2)}}{(q+2)!} \tilde{\Lambda}_{q,2} \right) \right\};$$

with  $q$  odd,

$$\sqrt{nh}h^{p+3} \left\{ e'_0 \tilde{\Gamma}_p^{-1} \left( \frac{m^{(p+2)}}{(p+2)!} \tilde{\Lambda}_{p,2} + \frac{m^{(p+3)}}{(p+3)!} \tilde{\Lambda}_{p,3} \right) - \rho^{-2} b^{q-(p+2)} \frac{m^{(q+1)}}{(q+1)!} e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} \tilde{\Lambda}_{q,1} \right\};$$

and finally at a boundary point, for any  $q$ ,

$$\eta_{\text{bc}} = \sqrt{nh}h^{p+2} \left\{ \frac{m^{(p+2)}}{(p+2)!} e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,2} - \rho^{-1} b^{q-(p+1)} \frac{m^{(q+1)}}{(q+1)!} e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} \tilde{\Lambda}_{q,1} \right\}.$$

## S.II.5 Main Result: Edgeworth Expansion

We now state our generic Edgeworth expansion, from whence the coverage probability expansion results follow immediately. We have opted to state separate results for undersmoothing, bias correction, and robust bias correction, rather than the unified statement of Theorem 3, for clarity. The unified structure is still present, and will be used in the proof of the result below, but is too cumbersome to use here. The Standard Normal distribution and density functions are  $\Phi(z)$  and  $\phi(z)$ , respectively.

**Theorem 5.** *Let Assumptions S.II.3.1, S.II.3.2, and S.II.3.3 hold, and assume  $nh/\log(n) \rightarrow \infty$ .*

(a) *If  $\eta_{\text{us}} \log(nh) \rightarrow 0$ , then for*

$$F_{\text{us}}(z) = \Phi(z) + \frac{1}{\sqrt{nh}} p_{1,\text{us}}(z) + \tilde{\eta}_{\text{us}} p_{3,\text{us}}(z) + \frac{1}{nh} q_{1,\text{us}}(z) + \tilde{\eta}_{\text{us}}^2 q_{2,\text{us}}(z) + \frac{\tilde{\eta}_{\text{us}}}{\sqrt{nh}} q_{3,\text{us}}(z),$$

*we have*

$$\sup_{z \in \mathbb{R}} |\mathbb{P}[T_{\text{us}} < z] - F_{\text{us}}(z)| = o\left((nh)^{-1} + (nh)^{-1/2} \tilde{\eta}_{\text{us}} + \tilde{\eta}_{\text{us}}^2\right).$$

(b) *If  $\eta_{\text{bc}} \log(nh) \rightarrow 0$  and  $\rho \rightarrow 0$ , then for*

$$F_{\text{bc}}(z) = \Phi(z) + \frac{1}{\sqrt{nh}} p_{1,\text{bc}}(z) + \tilde{\eta}_{\text{bc}} p_{3,\text{bc}}(z) + \frac{1}{nh} q_{1,\text{us}}(z) + \tilde{\eta}_{\text{bc}}^2 q_{2,\text{bc}}(z) + \frac{\tilde{\eta}_{\text{bc}}}{\sqrt{nh}} q_{3,\text{bc}}(z)$$

$$- \rho^{p+2}(\Omega_1 + \rho^{p+1}\Omega_2) \frac{\phi(z)}{2} z,$$

we have

$$\sup_{z \in \mathbb{R}} |\mathbb{P}[T_{\text{bc}} < z] - F_{\text{bc}}(z)| = o\left((nh)^{-1} + (nh)^{-1/2} \tilde{\eta}_{\text{bc}} + \tilde{\eta}_{\text{bc}}^2 + \rho^{1+2(p+1)}\right).$$

(c) If  $\eta_{\text{bc}} \log(nh) \rightarrow 0$  and  $\rho \rightarrow \bar{\rho} < \infty$ , then for

$$F_{\text{rbc}}(z) = \Phi(z) + \frac{1}{\sqrt{nh}} p_{1,\text{rbc}}(z) + \tilde{\eta}_{\text{bc}} p_{3,\text{rbc}}(z) + \frac{1}{nh} q_{1,\text{rbc}}(z) + \tilde{\eta}_{\text{bc}}^2 q_{2,\text{rbc}}(z) + \frac{\tilde{\eta}_{\text{bc}}}{\sqrt{nh}} q_{3,\text{rbc}}(z),$$

we have

$$\sup_{z \in \mathbb{R}} |\mathbb{P}[T_{\text{rbc}} < z] - F_{\text{rbc}}(z)| = o\left((nh)^{-1} + (nh)^{-1/2} \tilde{\eta}_{\text{bc}} + \tilde{\eta}_{\text{bc}}^2\right).$$

### S.II.5.1 Coverage Error for Undersmoothing

For undersmoothing estimators, we have the following result, which is valid for both interior and boundary points, with moments appropriately truncated if necessary. This result is the analogue of the robust bias correction corollary in the main text, and follows directly from the generic theorem there or Theorem 5 above. Exponents such as  $1 + 2(p+1)$  are intentionally not simplified to ease comparison to other results, particularly the density case.

The polynomials  $q_{1,\text{us}}$ ,  $q_{2,\text{us}}$ , and  $q_{3,\text{us}}$ , which do *not* have an argument, are defined in terms of those given in Section S.II.1.2 and used in Theorem 5, which *do* have an argument. Specifically, the polynomials in Section S.II.1.2 and Theorem 5 should be doubled, divided by the standard Normal density, and evaluated at the Normal quantile  $z_{\alpha/2}$ , that is,

$$q_{k,\text{us}} := \frac{2}{\phi(z)} q_{k,\text{us}}(z) \Big|_{z=z_{\alpha/2}}, \quad k = 1, 2, 3.$$

**Corollary 8** (Undersmoothing). *Let the conditions of Theorem 5(a) hold. Then*

$$\begin{aligned} \mathbb{P}[m \in I_{\text{us}}] = 1 - \alpha + \left\{ \frac{1}{nh} q_{1,\text{us}} + nh^{1+2(p+1)} (m^{(p+1)})^2 \left( e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} / (p+1)! \right)^2 q_{2,\text{us}} \right. \\ \left. + h^{p+1} (m^{(p+1)}) \left( e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} / (p+1)! \right) q_{3,\text{us}} \right\} \phi(z_{\alpha/2}) \{1 + o(1)\}. \end{aligned}$$

In particular, if  $h_{\text{us}}^* = H_{\text{us}}^* n^{-1/(1+(p+1))}$ , then  $\mathbb{P}[m \in I_{\text{us}}] = 1 - \alpha + O(n^{-(p+1)/(1+(p+1))})$ , where

$$H_{\text{us}}^* = \arg \min_H \left| H^{-1} q_{1,\text{us}} + H^{1+2(p+1)} (m^{(p+1)})^2 \left( e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} / (p+1)! \right)^2 q_{2,\text{us}} \right. \\ \left. + H^{p+1} (m^{(p+1)}) \left( e'_0 \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} / (p+1)! \right) q_{3,\text{us}} \right|.$$

## S.II.6 Proof of Main Result

We will first prove Theorem 5(a), as it is notationally simplest. From a technical and conceptual point of view, proving the remainder of Theorem 5 is identical, simply more involved notationally due to the additional complexity of the bias correction. Outlines of these proofs are found below.

### S.II.6.1 Proof of Theorem 5(a)

Let  $s = \sqrt{nh}$ .

Throughout this proof, we will generally omit the subscripts  $\text{us}$  and  $p$  when this causes no confusion. This entire proof focuses on the undersmoothing statistic,  $T_{\text{us}} = \hat{\sigma}_{\text{us}}^{-1} s(\hat{m} - m)$ , and since bias correction is not involved at all, the associated constructions such as  $\Gamma_q$ ,  $W_q$ , etc, do not appear, and hence there is no need to carry the additional notation to distinguish  $W_p$  from  $W_q$ , or  $\hat{\sigma}_{\text{us}}$  from  $\hat{\sigma}_{\text{rbc}}$ , for example, and we will simply write  $\Gamma$  for  $\Gamma_p$ ,  $W$  for  $W_p$ ,  $\hat{\sigma}$  for  $\hat{\sigma}_{\text{us}}$ , etc.

Our goal is to expand  $\mathbb{P}[T_{\text{us}} < z]$ , where  $T_{\text{us}} = \hat{\sigma}^{-1} s(\hat{m} - m)$ . The proof proceeds by identifying a smooth function  $\tilde{T} = \tilde{T}(z)$  such that, for the random variable  $Z_{\text{us}} := Z_{\text{us}}(u)$  that obeys Cramér's condition (Assumption S.II.3.3),  $\tilde{T}(\mathbb{E}[Z_{\text{us}}]) = 0$  and

$$\mathbb{P}[T_{\text{us}} < z] = \mathbb{P}[\tilde{T}(\bar{Z}_{\text{us}}) < \tilde{z}] + o(s^{-2} + s^{-1}\eta + \eta^2), \quad (23)$$

where  $\bar{Z} = \sum_{i=1}^n Z_i/n$  and  $\tilde{z}$  is a known, nonrandom quantity that depends on the original quantile  $z$  and the remainder  $T_{\text{us}} - \tilde{T}$  (see Remark 9). An Edgeworth expansion for  $\tilde{T}$  holds under Assumption S.II.3.3, and a Taylor expansion of this function around  $\tilde{z}$  yields the final result. As in the density case,  $\tilde{z}$  will capture the bias terms of  $T_{\text{us}}$ : in that case  $\tilde{z} = z - \eta/\tilde{\sigma}$ , but here bias is present in both the numerator and the Studentization.

To begin, define the notation  $\check{R} = [r_p(X_1 - x), \dots, r_p(X_n - x)]'$  and  $M = [m(X_1), \dots, m(X_n)]'$ , and use this to split  $T$  into variance and bias terms, as follows:

$$T = \hat{\sigma}^{-1} s e'_0 \Gamma^{-1} R' W (Y - M)/n + \hat{\sigma}^{-1} s e'_0 \Gamma^{-1} R' W (M - \check{R}\beta)/n.$$

We use this decomposition to rewrite  $\mathbb{P}[T_{\text{us}} < z]$  as

$$\begin{aligned}
\mathbb{P}[T_{\text{us}} < z] &= \mathbb{P}[T_{\text{us}} - \tilde{\sigma}^{-1}\eta < z - \tilde{\sigma}^{-1}\eta] \\
&= \mathbb{P}\left[\left\{\hat{\sigma}^{-1}se'_0\Gamma^{-1}R'W(Y - M)/n + \hat{\sigma}^{-1}se'_0\Gamma^{-1}R'W(M - \check{R}\beta)/n - \tilde{\sigma}^{-1}\eta\right\} < z - \tilde{\sigma}^{-1}\eta\right] \\
&= \mathbb{P}\left[\left\{\tilde{\sigma}^{-1}se'_0\Gamma^{-1}R'W(Y - M)/n \right. \right. \\
&\quad \left. \left. + \tilde{\sigma}^{-1}se'_0\tilde{\Gamma}^{-1}R'W(M - \check{R}\beta)/n - \tilde{\sigma}^{-1}\eta \right. \right. \\
&\quad \left. \left. + \tilde{\sigma}^{-1}se'_0\left(\Gamma^{-1} - \tilde{\Gamma}^{-1}\right)R'W(M - \check{R}\beta)/n \right. \right. \\
&\quad \left. \left. + (\hat{\sigma}^{-1} - \tilde{\sigma}^{-1})se'_0\Gamma^{-1}R'W(Y - M)/n \right. \right. \\
&\quad \left. \left. + (\hat{\sigma}^{-1} - \tilde{\sigma}^{-1})se'_0\Gamma^{-1}R'W(M - \check{R}\beta)/n\right\} < z - \tilde{\sigma}^{-1}\eta\right]. \tag{24}
\end{aligned}$$

The first three lines in the last equality obey the desired properties of  $\tilde{T}$  by the orthogonality of  $\varepsilon_i$ , the definition of  $\eta_{\text{us}}$  in Eqn. (22) as  $\mathbb{E}\left[se'_0\tilde{\Gamma}^{-1}R'W(M - \check{R}\beta)/n\right]$ , and the fact that  $\Gamma^{-1} - \tilde{\Gamma}^{-1} = \tilde{\Gamma}^{-1}(\tilde{\Gamma} - \Gamma)\Gamma^{-1}$ . For the final two (which are  $T_{\text{us}} - \tilde{\sigma}^{-1}s(\hat{m} - m) = \hat{\sigma}^{-1} - \tilde{\sigma}^{-1}s(\hat{m} - m)$ ), we must expand the difference  $\hat{\sigma}^{-1} - \tilde{\sigma}^{-1}$ . Accounting for the resulting terms will constitute the bulk of the remainder of the proof, as well as complete the construction of  $\tilde{z}$  and the remainder terms of Eqn. (23).<sup>2</sup>

To begin, with  $\tilde{\sigma}^2 = e'_0\tilde{\Gamma}^{-1}\tilde{\Psi}\tilde{\Gamma}^{-1}e_0$  defined in Section S.II.1.2,

$$\frac{1}{\hat{\sigma}} = \frac{1}{\tilde{\sigma}} \left( \frac{\hat{\sigma}^2}{\tilde{\sigma}^2} \right)^{-1/2} = \frac{1}{\tilde{\sigma}} \left( 1 + \frac{\hat{\sigma}^2 - \tilde{\sigma}^2}{\tilde{\sigma}^2} \right)^{-1/2},$$

and hence a Taylor expansion gives

$$\frac{1}{\hat{\sigma}} = \frac{1}{\tilde{\sigma}} \left[ 1 - \frac{1}{2} \frac{\hat{\sigma}^2 - \tilde{\sigma}^2}{\tilde{\sigma}^2} + \frac{3}{8} \left( \frac{\hat{\sigma}^2 - \tilde{\sigma}^2}{\tilde{\sigma}^2} \right)^2 - \frac{1}{3!} \frac{15}{8} \left( \frac{\hat{\sigma}^2 - \tilde{\sigma}^2}{\tilde{\sigma}^2} \right)^3 \frac{\tilde{\sigma}^7}{\tilde{\sigma}^7} \right],$$

for a point  $\tilde{\sigma}^2 \in [\tilde{\sigma}^2, \hat{\sigma}^2]$ , and so

$$\hat{\sigma}^{-1} - \tilde{\sigma}^{-1} = -\frac{1}{2} \frac{\hat{\sigma}^2 - \tilde{\sigma}^2}{\tilde{\sigma}^3} + \frac{3}{8} \frac{(\hat{\sigma}^2 - \tilde{\sigma}^2)^2}{\tilde{\sigma}^5} - \frac{5}{16} \frac{(\hat{\sigma}^2 - \tilde{\sigma}^2)^3}{\tilde{\sigma}^7}. \tag{25}$$

We thus focus on  $\hat{\sigma}^2 - \tilde{\sigma}^2$ . Recall the definition of  $\tilde{\Psi} = hR'W\Sigma WR/n$ . Then define the two

---

<sup>2</sup>Technically, to obtain a  $\tilde{T}$  with the desired properties, one need not expand  $\hat{\sigma}^{-1} - \tilde{\sigma}^{-1}$  for the variance term: that is, in Eqn. (24),  $\tilde{\sigma}^{-1}se'_0\Gamma^{-1}R'W(Y - M)/n$  and  $(\hat{\sigma}^{-1} - \tilde{\sigma}^{-1})se'_0\Gamma^{-1}R'W(Y - M)/n$  may be collapsed. This requires strengthening Cramér's condition (see Remark 9), and since  $\hat{\sigma}^{-1} - \tilde{\sigma}^{-1}$  must be accounted for in the final bias term,  $(\hat{\sigma}^{-1} - \tilde{\sigma}^{-1})se'_0\Gamma^{-1}R'W(M - \check{R}\beta)/n$ , there is little reason not to do both terms.

terms  $A_1$  and  $A_2$  through the following:

$$\hat{\sigma}^2 - \tilde{\sigma}^2 = e_0' \Gamma^{-1} \left( \hat{\Psi} - \check{\Psi} \right) \Gamma^{-1} e_0 + \left( e_0' \Gamma^{-1} \check{\Psi} \Gamma^{-1} e_0 - e_0' \tilde{\Gamma}^{-1} \tilde{\Psi} \tilde{\Gamma}^{-1} e_0 \right) =: A_1 + A_2. \quad (26)$$

For  $A_1$ , recall that  $\hat{\varepsilon}_i = y_i - r_p(X_i - x)' \hat{\beta}_p$  and so

$$\begin{aligned} \hat{\Psi} - \check{\Psi} &= \frac{1}{nh} \sum_{i=1}^n (K r_p r_p')(X_{h,i}) \{ \hat{\varepsilon}_i^2 - v(X_i) \} \\ &= \frac{1}{nh} \sum_{i=1}^n (K r_p r_p')(X_{h,i}) \left\{ \left( y_i - r_p(X_i - x)' \hat{\beta}_p \right)^2 - v(X_i) \right\} \\ &= \frac{1}{nh} \sum_{i=1}^n (K r_p r_p')(X_{h,i}) \left\{ \left( \varepsilon_i + [m(X_i) - r_p(X_i - x)' \beta_p] + r_p(X_i - x)' [\beta_p - \hat{\beta}_p] \right)^2 - v(X_i) \right\} \\ &=: A_{1,1} + A_{1,2} + A_{1,3} + A_{1,4} + A_{1,5} + A_{1,6} + A_{1,7} + A_{1,8}, \end{aligned} \quad (27)$$

where

$$A_{1,1} = \frac{1}{nh} \sum_{i=1}^n (K r_p r_p')(X_{h,i}) \{ \varepsilon_i^2 - v(X_i) \},$$

is due to the approximation of the (average over the) conditional variance by the squared residuals (i.e.  $A_{1,1}$  is the sole remainder that would arise if the true residuals were known and used in place of  $\hat{\varepsilon}_i^2$ ), and, using  $r_p(X_i - x)' \hat{\beta} = r_p(X_i - x)' H_p \Gamma^{-1} R' W Y / n = r_p(X_{h,i})' \Gamma^{-1} R' W Y / n$ , the terms  $A_{1,k}$ ,  $k = 2, 3, \dots, 8$  are:

$$\begin{aligned} A_{1,2} &= \frac{1}{nh} \sum_{i=1}^n (K r_p r_p')(X_{h,i}) \{ 2\varepsilon_i [m(X_i) - r_p(X_i - x)' \beta_p] \}, \\ A_{1,3} &= \frac{1}{nh} \sum_{i=1}^n (K r_p r_p')(X_{h,i}) \{ -2\varepsilon_i r_p(X_{h,i})' \} \Gamma^{-1} R' W (Y - \check{R}\beta) / n, \\ A_{1,4} &= \frac{1}{nh} \sum_{i=1}^n (K r_p r_p')(X_{h,i}) \{ -2[m(X_i) - r_p(X_i - x)' \beta_p] r_p(X_{h,i})' \} \Gamma^{-1} R' W (Y - M) / n, \\ A_{1,5} &= \frac{1}{nh} \sum_{i=1}^n (K r_p r_p' r_p')(X_{h,i}) \Gamma^{-1} R' W (Y - M) / n \left[ (Y - M)' / n + 2(M - \check{R}\beta) / n \right] W R \Gamma^{-1} r_p(X_{h,i}), \\ A_{1,6} &= \frac{1}{nh} \sum_{i=1}^n (K r_p r_p')(X_{h,i}) [m(X_i) - r_p(X_i - x)' \beta_p]^2, \\ A_{1,7} &= \frac{1}{nh} \sum_{i=1}^n (K r_p r_p' r_p')(X_{h,i}) \{ -2[m(X_i) - r_p(X_i - x)' \beta_p] \} \Gamma^{-1} R' W (M - \check{R}\beta) / n, \end{aligned}$$



and

$$A_{1,8} = \frac{1}{nh} \sum_{i=1}^n (K r_p r_p' r_p') (X_{h,i}) \Gamma^{-1} [R' W (M - \check{R}\beta)/n] [(M - \check{R}\beta)' / n W R] \Gamma^{-1} r_p (X_{h,i}).$$

With this notation, we can write  $A_1 = e_0' \Gamma^{-1} (\hat{\Psi} - \check{\Psi}) \Gamma^{-1} e_0 = e_0' \Gamma^{-1} (\sum_{k=1}^8 A_{1,k}) \Gamma^{-1} e_0$ . The terms  $A_{1,1}$  to  $A_{1,5}$  will be incorporated into  $\tilde{T}$ : notice that these terms obey  $A_{1,k} = A_{1,k}(\bar{Z}_{\text{us}})$  and  $A_{1,k}(\mathbb{E}[Z_{\text{us}}]) = 0$ , and hence these properties will be inherited in the final two lines of Eqn. (24). However,  $A_{1,6}$ ,  $A_{1,7}$ , and  $A_{1,8}$  do not have these properties, and will thus be incorporated into  $\tilde{z}$  and the remainder. Details are below.

Turning to  $A_2$  in Eqn. (26), using the identity  $\Gamma^{-1} - \tilde{\Gamma}^{-1} = \tilde{\Gamma}^{-1} (\tilde{\Gamma} - \Gamma) \Gamma^{-1}$  and that  $\Gamma$  and  $\Psi$  are symmetric, we find that

$$\begin{aligned} A_2 &= e_0' \Gamma^{-1} \check{\Psi} \Gamma^{-1} e_0 - e_0' \tilde{\Gamma}^{-1} \tilde{\Psi} \tilde{\Gamma}^{-1} e_0 \\ &= e_0' \Gamma^{-1} (\check{\Psi} - \tilde{\Psi}) \Gamma^{-1} e_0 + e_0' (\Gamma^{-1} - \tilde{\Gamma}^{-1}) \tilde{\Psi} \tilde{\Gamma}^{-1} e_0 + e_0' (\Gamma^{-1} - \tilde{\Gamma}^{-1}) \tilde{\Psi} \tilde{\Gamma}^{-1} e_0 \\ &= e_0' \Gamma^{-1} (\check{\Psi} - \tilde{\Psi}) \Gamma^{-1} e_0 - e_0' \tilde{\Gamma}^{-1} (\Gamma - \tilde{\Gamma}) \Gamma^{-1} \tilde{\Psi} (\Gamma^{-1} + \tilde{\Gamma}^{-1}) e_0. \end{aligned}$$

All of these terms obey the required properties of  $\tilde{T}$ .

We now collect the terms from expanding  $\hat{\sigma}^{-1} - \tilde{\sigma}^{-1}$  and return to Eqn. (24). Plugging the terms  $A_{1,1}-A_{1,8}$  and  $A_2$  into the Taylor expansion in Eqn. (25), by way of Eqn. (26), and collecting terms appropriately (i.e. those that belong in  $\tilde{T}$  as described above), we have the following, which picks up from Eqn. (24) and is a precursor to Eqn. (23):

$$\mathbb{P}[T_{\text{us}} < z] = \mathbb{P} \left[ \tilde{T}(\bar{Z}_{\text{us}}) + U < \tilde{z} \right]. \quad (28)$$

In this statement, we have made the following constructions:

$$\begin{aligned} \tilde{T} &= \tilde{\sigma}^{-1} s e_0' \Gamma^{-1} R' W (Y - M) / n \\ &\quad + \tilde{\sigma}^{-1} s e_0' \tilde{\Gamma}^{-1} R' W (M - \check{R}\beta) / n - \tilde{\sigma}^{-1} \eta \\ &\quad + \tilde{\sigma}^{-1} s e_0' (\Gamma^{-1} - \tilde{\Gamma}^{-1}) R' W (M - \check{R}\beta) / n \\ &\quad + \left\{ -\frac{1}{2\tilde{\sigma}^3} \left[ e_0' \Gamma^{-1} \left( \sum_{k=1}^5 A_{1,k} \right) \Gamma^{-1} e_0 + A_2 \right] + \frac{3}{8\tilde{\sigma}^5} \left[ e_0' \Gamma^{-1} A_{1,1} \Gamma^{-1} e_0 + A_2 \right]^2 \right\} \\ &\quad \times \left\{ s e_0' \Gamma^{-1} R' W (Y - M) / n + s e_0' \Gamma^{-1} R' W (M - \check{R}\beta) / n \right\}, \end{aligned}$$

$$\begin{aligned}
U = & \left\{ -\frac{1}{2\tilde{\sigma}^3} e'_0 \Gamma^{-1} (A_{1,6} + A_{1,7} + A_{1,8}) \Gamma^{-1} e_0 + \frac{3}{8\tilde{\sigma}^5} \left[ e'_0 \Gamma^{-1} \left( \sum_{k=2}^8 A_{1,k} \right) \Gamma^{-1} e_0 \right]^2 - \frac{5}{16} \frac{(\hat{\sigma}^2 - \tilde{\sigma}^2)^3}{\tilde{\sigma}^7} \right\} \\
& \times \left\{ s e'_0 \Gamma^{-1} R' W (Y - M) / n + s e'_0 \Gamma^{-1} R' W (M - \tilde{R} \beta) / n \right\} \\
& - \left\{ -\frac{1}{2\tilde{\sigma}^3} e'_0 \tilde{\Gamma}^{-1} \left( \tilde{A}_{1,6} + \tilde{A}_{1,7} + \tilde{A}_{1,8} \right) \tilde{\Gamma}^{-1} e_0 \right\} \eta,
\end{aligned}$$

and

$$\tilde{z} = z - \left\{ \tilde{\sigma}^{-1} - \frac{1}{2\tilde{\sigma}^3} e'_0 \tilde{\Gamma}^{-1} \left( \tilde{A}_{1,6} + \tilde{A}_{1,7} + \tilde{A}_{1,8} \right) \tilde{\Gamma}^{-1} e_0 \right\} \eta.$$

In  $U$  and  $\tilde{z}$ , each  $\tilde{A}_{1,k}$  is  $A_{1,k}$  where all elements have been replaced by their respective fixed- $n$  expected values, that is,

$$\begin{aligned}
\tilde{A}_{1,6} &= \mathbb{E}[A_{1,6}] = \mathbb{E} \left[ h^{-1} (K r_p r'_p) (X_{h,i}) [m(X_i) - r_p(X_i - x)' \beta_p]^2 \right], \\
\tilde{A}_{1,7} &= -2 \mathbb{E} \left[ h^{-1} (K r_p r'_p r'_p) (X_{h,i}) [m(X_i) - r_p(X_i - x)' \beta_p] \right] \\
&\quad \times \tilde{\Gamma}^{-1} \mathbb{E} \left[ h^{-1} (K r_p) (X_{h,j}) [m(X_j) - r_p(X_j - x)' \beta_p] \right],
\end{aligned}$$

and

$$\tilde{A}_{1,8} = \mathbb{E} \left[ h^{-1} (K r_p r'_p) (X_{h,i}) \mathbb{E} \left[ h^{-1} r_p (X_{h,i}) \tilde{\Gamma}^{-1} (K r_p) (X_{h,j}) [m(X_j) - r_p(X_j - x)' \beta_p] \middle| X_i \right]^2 \right].$$

The next step in the proof is to show that, for  $r_* = \max\{s^{-2}, \eta^2, h^{p+1}\}$  (i.e., the slowest decaying), it holds that

$$\frac{1}{r_*} \mathbb{P}[|U| > r_n] \rightarrow 0, \quad \text{for some } r_n = o(r_*). \quad (29)$$

This result is established by Lemma 7 in Section S.II.6.3 below. This, together with Eqn. (28), implies Eqn. (23).

Under Assumption S.II.3.3, an Edgeworth expansion holds for  $\tilde{T}$  up to  $o(s^{-2} + s^{-1}\eta + \eta^2)$ . Thus, for a smooth function  $G(z)$ , we have  $\mathbb{P}[\tilde{T} < z] = G(z) + o(s^{-2} + s^{-1}\eta + \eta^2)$ . Therefore, a Taylor expansion gives

$$\mathbb{P}[\tilde{T} < \tilde{z}] = G(z) - G^{(1)}(z) \left\{ \tilde{\sigma}^{-1} - \frac{1}{2\tilde{\sigma}^3} e'_0 \Gamma^{-1} \left( \tilde{A}_{1,6} + \tilde{A}_{1,7} + \tilde{A}_{1,8} \right) \Gamma^{-1} e_0 \right\} + o(s^{-2} + s^{-1}\eta + \eta^2),$$

which together with Eqn. (23) establishes the validity of the Edgeworth expansion. The terms of the expansion are computed in Section S.II.6.4 below.  $\square$

### S.II.6.2 Proof of Theorem 5(b) & (c)

To prove parts (b) and (c) of Theorem 5 the same steps are required, and so we will not pursue all the details here. Indeed, the same expansions are performed and the same bounds computed on objects which are conceptually similar, only taking into account the bias correction (in the numerator for (b), and also in the denominator for (c)). The bias correction will result in essentially two changes: first, many more terms like  $\Gamma - \tilde{\Gamma}$  appear, and second, the bias expressions and rates change. To illustrate, we will list several key points where these changes manifest. This list is not exhaustive, but it will show that the same methods used above still apply.

First, for the numerator of  $T_{bc}$  and  $T_{rbc}$ , recall that the estimator  $\hat{m}$  is

$$\hat{m} = \left\{ e'_0 \Gamma_p^{-1} R'_p W_p \right\} Y/n,$$

while the bias corrected estimator is

$$\hat{m} - \hat{B}_m = \left\{ e'_0 \Gamma_p^{-1} (R'_p W_p - \rho^{p+1} \Lambda_{p,1} e'_{p+1} \Gamma_q^{-1} R'_q W_q) \right\} Y/n.$$

Comparing these two expressions, it can be seen that the terms in the proof above that involve  $\Gamma_p - \tilde{\Gamma}_p$  will now additionally involve  $\Gamma_q - \tilde{\Gamma}_q$  and  $\Lambda_{p,1} - \tilde{\Lambda}_{p,1}$ , whereas those that with  $e'_0 \tilde{\Gamma}_p^{-1} R'_p W_p$  will now have  $e'_0 \tilde{\Gamma}_p^{-1} (R'_p W_p - \rho^{p+1} \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} R'_q W_q)$  instead. To give a concrete example, consider the third line of Eqn. (24),

$$\tilde{\sigma}_{us}^{-1} s e'_0 \left( \Gamma_p^{-1} - \tilde{\Gamma}_p^{-1} \right) R'_p W_p (M - \check{R}_p \beta_p)/n,$$

which becomes a piece of the function  $\tilde{T}$ . For part (b) Theorem 5, treating  $T_{bc}$ , this will become

$$\begin{aligned} & \tilde{\sigma}_{us}^{-1} s e'_0 \left( \Gamma_p^{-1} - \tilde{\Gamma}_p^{-1} \right) R'_p W_p (M - \check{R}_{p+1} \beta_{p+1})/n \\ & - s e'_0 \rho^{p+1} \left( \Gamma_p^{-1} \Lambda_{p,1} e'_{p+1} \Gamma_q^{-1} - \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} \right) R'_q W_q (M - \check{R}_q \beta_q)/n, \end{aligned}$$

and part (c) will have the same but with  $\tilde{\sigma}_{rbc}^{-1}$ . Then, since

$$\begin{aligned} \Gamma_p^{-1} \Lambda_{p,1} e'_{p+1} \Gamma_q^{-1} - \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e'_{p+1} \tilde{\Gamma}_q^{-1} &= \left( \Gamma_p^{-1} - \tilde{\Gamma}_p^{-1} \right) \Lambda_{p,1} e'_{p+1} \Gamma_q^{-1} \\ &+ \tilde{\Gamma}_p^{-1} \left( \Lambda_{p,1} - \tilde{\Lambda}_{p,1} \right) e'_{p+1} \Gamma_q^{-1} + \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1} e'_{p+1} \left( \Gamma_q^{-1} - \tilde{\Gamma}_q^{-1} \right), \end{aligned}$$

this term is handled identically, since the appropriate Cramér's condition is assumed.

Consider now the denominator of the Studentized statistics. For part (b), there is no change as  $\hat{\sigma}_{\text{us}}^2$  is still used, and so the terms involving  $A_{1,k}$  and  $A_2$  will be identical. However, for  $T_{\text{rbc}}$ , we must account for changes of the above form, but also that the residuals are estimated with the degree  $q$  fit:  $\hat{\varepsilon}_i = y_i - r_q(X_i - x)' \hat{\beta}_q$  instead of degree  $p$ . With these changes in mind, the analogue of Eqn. (26) will be

$$\hat{\sigma}_{\text{rbc}}^2 - \tilde{\sigma}_{\text{rbc}}^2 = e_0' \Gamma_p^{-1} \left( \hat{\Psi}_q - \check{\Psi}_q \right) \Gamma_p^{-1} e_0 + \left( e_0' \Gamma_p^{-1} \check{\Psi}_q \Gamma_p^{-1} e_0 - e_0' \tilde{\Gamma}_p^{-1} \tilde{\Psi}_q \tilde{\Gamma}_p^{-1} e_0 \right). \quad (30)$$

The second term will proceed as above, though  $\hat{\Psi}_p - \check{\Psi}_p$  will be replaced by

$$\hat{\Psi}_q - \check{\Psi}_q = \frac{1}{nh} \sum_{i=1}^n \left\{ \tilde{\ell}_{\text{bc}}^0(X_i) \tilde{\ell}_{\text{bc}}^0(X_i)' v(X_i) - \mathbb{E} \left[ \tilde{\ell}_{\text{bc}}^0(X_i) \tilde{\ell}_{\text{bc}}^0(X_i)' v(X_i) \right] \right\},$$

where  $\tilde{\ell}_{\text{bc}}^0(X_i) = (Kr_p)(X_{h,i}) - \rho^{p+2} \tilde{\Lambda}_{p,1} \tilde{\Gamma}_q^{-1} (Lr_p)(\rho X_{h,i})$  (cf. Section S.II.1.2, the function  $\ell_{\text{bc}}^0$  therein is  $\ell_{\text{bc}}^0(X_i) = e_0' \tilde{\Gamma}_p^{-1} \tilde{\ell}_{\text{bc}}^0(X_i)$ ). To use similar notation,

$$\hat{\Psi}_p - \check{\Psi}_p = \frac{1}{nh} \sum_{i=1}^n \left\{ \tilde{\ell}_{\text{us}}^0(X_i) \tilde{\ell}_{\text{us}}^0(X_i)' v(X_i) - \mathbb{E} \left[ \tilde{\ell}_{\text{us}}^0(X_i) \tilde{\ell}_{\text{us}}^0(X_i)' v(X_i) \right] \right\}.$$

Then, expanding  $\tilde{\ell}_{\text{bc}}^0(X_i)$  shows that  $\hat{\Psi}_q - \check{\Psi}_q$  is equal to

$$\begin{aligned} \check{\Psi}_p - \tilde{\Psi}_p + \rho^{(p+1)+1} \tilde{\Lambda}_{p,1} \tilde{\Gamma}_q^{-1} \frac{1}{nb} \sum_{i=1}^n \left\{ (Lr_q r_q')(X_{b,i}) v(X_i) - \mathbb{E} \left[ (Lr_q r_q')(X_{b,i}) v(X_i) \right] \right\} \tilde{\Gamma}_q^{-1} \tilde{\Lambda}_{p,1} \\ - \rho^{(p+1)+1} 2 \frac{1}{nh} \sum_{i=1}^n \left\{ (Kr_p)(X_{h,i}) (Lr_q')(\rho X_{h,i}) v(X_i) - \mathbb{E} \left[ (Kr_p)(X_{h,i}) (Lr_q')(\rho X_{h,i}) v(X_i) \right] \right\} \tilde{\Gamma}_q^{-1} \tilde{\Lambda}_{p,1}, \end{aligned}$$

and since all these terms still obey the appropriate Cramér's condition, the same steps apply.

The first term of Eqn. (30) will also follow by the same method as in the prior proof, but more care must be taken as many more terms will be present because  $\hat{\Psi}_q - \check{\Psi}_q$  consists of the following three terms, representing the variance of  $\hat{m}$ , the variance of  $\hat{B}_m$ , and their covariance, respectively:

$$\begin{aligned} \hat{\Psi}_q - \check{\Psi}_q &= h R_p' W_p \left( \hat{\Sigma}_q - \Sigma \right) W_p R_p / n \\ &\quad + h \rho^{2(p+1)} \Lambda_{p,1} \Gamma_q^{-1} \left( R_q' W_q \hat{\Sigma}_q W_q R_q \right) \Gamma_q^{-1} \Lambda_{p,1}' / n - h \rho^{2(p+1)} \tilde{\Lambda}_{p,1} \tilde{\Gamma}_q^{-1} \left( R_q' W_q \Sigma W_q R_q \right) \tilde{\Gamma}_q^{-1} \tilde{\Lambda}_{p,1}' / n \\ &\quad - 2h \rho^{p+1} R_p' W_p \left( \hat{\Sigma}_q W_q R_q \Gamma_p^{-1} \Lambda_{p,1}' \Gamma - \Sigma W_q R_q \tilde{\Gamma}_p^{-1} \tilde{\Lambda}_{p,1}' \right) / n. \end{aligned}$$

The first of these three is as in the prior proof, and yields the same  $A_{1,1}$ – $A_{1,8}$ , only with the

bias of a  $q$ -degree fit:  $m(X_i) - r_q(X_i - x)' \beta_q$ . If we define

$$\check{\Psi}_q := \frac{1}{nb} \sum_{i=1}^n (L^2 r_q r_q')(X_{b,i}) v(X_i)$$

then the second term of  $\hat{\Psi}_q - \check{\Psi}_q$  is equal to

$$\begin{aligned} & \rho^{1+2(p+1)} \Lambda_{p,1} \Gamma_q^{-1} \left\{ \frac{1}{nb} \sum_{i=1}^n (L^2 r_q r_q')(X_{b,i}) \{ \hat{\varepsilon}_i^2 - v(X_i) \} \right\} \Gamma_q^{-1} \Lambda_{p,1} \\ & + \rho^{1+2(p+1)} \left( \Lambda_{p,1} - \tilde{\Lambda}_{p,1} \right) \Gamma_q^{-1} \check{\Psi}_q \Gamma_q^{-1} \Lambda_{p,1} \\ & + \rho^{1+2(p+1)} \tilde{\Lambda}_{p,1} \left( \Gamma_q^{-1} - \tilde{\Gamma}_q^{-1} \right) \check{\Psi}_q \Gamma_q^{-1} \Lambda_{p,1} \\ & + \rho^{1+2(p+1)} \tilde{\Lambda}_{p,1} \tilde{\Gamma}_q^{-1} \check{\Psi}_q \left( \Gamma_q^{-1} - \tilde{\Gamma}_q^{-1} \right) \Lambda_{p,1} \\ & + \rho^{1+2(p+1)} \tilde{\Lambda}_{p,1} \tilde{\Gamma}_q^{-1} \check{\Psi}_q \tilde{\Gamma}_q^{-1} \left( \Lambda_{p,1} - \tilde{\Lambda}_{p,1} \right). \end{aligned}$$

The first of these terms will also give rise to versions of  $A_{1,1}-A_{1,8}$ , only with the bias of a  $q$ -degree fit and changing  $K$  to  $L$ ,  $p$  to  $q$ ,  $h$  to  $b$ , etc, and will thus be treated exactly as above. The rest of these are incorporated into  $\tilde{T}_{\text{rbc}}$ , similar to how  $A_2$  is treated, because Cramér's condition is satisfied. The third and final piece of  $\hat{\Psi}_q - \check{\Psi}_q$  is equal to

$$\begin{aligned} & - 2\rho^{1+(p+1)} \left\{ \frac{1}{nh} \sum_{i=1}^n (K r_p)(X_{h,i}) (L r_q')(X_{h,i} \rho) \{ \hat{\varepsilon}_i^2 - v(X_i) \} \right\} \Gamma_q^{-1} \Lambda'_{p,1} \\ & - 2\rho^{1+(p+1)} \check{\Psi}_q \left( \Gamma_q^{-1} - \tilde{\Gamma}_q^{-1} \right) \Lambda'_{p,1} \\ & - 2\rho^{1+(p+1)} \check{\Psi}_q \tilde{\Gamma}_q^{-1} \left( \Lambda_{p,1} - \tilde{\Lambda}_{p,1} \right), \end{aligned}$$

and thus is entirely analogous, with yet another version of  $A_{1,1}-A_{1,8}$  defined for the remainder in the first line, and the second two easily incorporated into  $\tilde{T}_{\text{rbc}}$ .

From these arguments, it is clear that the analogue of Lemma 7 will hold for these cases as well: the same fundamental pieces are involved, and thus the same arguments will apply, just as above.

### S.II.6.3 Lemmas

Our proof of Theorem 5 relies on the following lemmas. The first gives generic results used to derive rate bounds on the probability of deviations of the necessary terms. Some such results are collected in Lemma 5. Lemma 7 shows how to use the previous results to establish

negligibility of the remainder terms required for Eqn. (29).

As above, we will generally omit the details required for Theorem 5 parts (b) and (c), to save space. These are entirely analogous, as can be seen from the steps in Lemma 5. Indeed, the first results are stated in terms of the kernel  $K$  and bandwidth  $h$ , but continue to hold for  $L$  and  $b$  under the obvious substitutions and appropriate assumptions.

Throughout proofs  $C$  shall be a generic conformable constant that may take different values in different places. If more than one constant is needed,  $C_1, C_2, \dots$ , will be used.

**Lemma 4.** *Let the conditions of Theorem 5 hold and let  $g(\cdot)$  and  $t(\cdot)$  be continuous scalar functions.*

(a) *For some  $\delta > 0$ ,*

$$s^2 \mathbb{P} \left[ \left| s^{-2} \sum_{i=1}^n \{ (Kt)(X_{h,i})g(X_i) - \mathbb{E}[(Kt)(X_{h,i})g(X_i)] \} \right| > \delta s^{-1} \log(s)^{1/2} \right] \rightarrow 0.$$

(b) *For some  $\delta > 0$ ,*

$$s^2 \mathbb{P} \left[ \left| s^{-1} \sum_{i=1}^n \{ (Kt)(X_{h,i})g(X_i)\varepsilon_i \} \right| > \delta \log(s)^{1/2} \right] \rightarrow 0.$$

*The same holds with  $\varepsilon_i^2 - v(X_i)$  in place of  $\varepsilon_i$ , since it is conditionally mean zero and has more than four moments.*

(c) *For any  $\delta > 0$ , an integer  $k$ , and any  $\gamma > 0$ ,*

$$\frac{1}{h^{p+1}} \mathbb{P} \left[ \left| s^{-2} \sum_{i=1}^n (Kt)(X_{h,i})g(X_i) [m(X_i) - r_p(X_i - x)' \beta_p]^k \right| > \delta h^{(k-1)(p+1)} \log(s)^\gamma \right] \rightarrow 0.$$

(d) *For any  $\delta > 0$  and any  $\gamma > 0$ ,*

$$s^2 \mathbb{P} \left[ \left| s^{-2} \sum_{i=1}^n (Kt)(X_{h,i})g(X_i)\varepsilon_i [m(X_i) - r_p(X_i - x)' \beta_p] \right| > \delta h^{p+1} \log(s)^\gamma \right] \rightarrow 0.$$

(e) *For any  $\delta > 0$ , an integer  $k$ , and any  $\gamma > 0$ ,*

$$s^2 \mathbb{P} \left[ \left| s^{-2} \sum_{i=1}^n \left\{ (Kt)(X_{h,i})g(X_i)(m(X_i) - r_p(X_i - x)' \beta_p)^k - \mathbb{E}[(Kt)(X_{h,i})g(X_i)(m(X_i) - r_p(X_i - x)' \beta_p)^k] \right\} \right| > \delta h^{k(p+1)} \log(s)^\gamma \right] \rightarrow 0.$$

*Proof of Lemma 4(a).* Because the kernel function has compact support and  $t$  and  $g$  are continuous, we have

$$|(Kt)(X_{h,i})g(X_i) - \mathbb{E}[(Kt)(X_{h,i})g(X_i)]| < C_1.$$

Further, by a change of variables and using the assumptions on  $f$ ,  $g$  and  $t$ :

$$\begin{aligned} \mathbb{V}[(Kt)(X_{h,i})g(X_i)] &\leq \mathbb{E}[(Kt)(X_{h,i})^2g(X_i)^2] = \int f(X_i)(Kt)(X_{h,i})^2g(X_i)^2dX_i \\ &= h \int f(x - uh)g(x - uh)(Kt)(u)^2du \leq C_2h. \end{aligned}$$

Therefore, by Bernstein's inequality

$$\begin{aligned} s^2\mathbb{P}\left[\left|\frac{1}{s^2}\sum_{i=1}^n\{(Kt)(X_{h,i})g(X_i) - \mathbb{E}[(Kt)(X_{h,i})g(X_i)]\}\right| > \delta s^{-1}\log(s)^{1/2}\right] \\ \leq 2s^2\exp\left\{-\frac{(s^4)(\delta s^{-1}\log(s)^{1/2})^2/2}{C_2s^2 + C_1s^2\delta s^{-1}\log(s)^{1/2}/3}\right\} \\ = 2\exp\{2\log(s)\}\exp\left\{-\frac{\delta^2\log(s)/2}{C_2 + C_1\delta s^{-1}\log(s)^{1/2}/3}\right\} \\ = 2\exp\left\{\log(s)\left[2 - \frac{\delta^2/2}{C_2 + C_1\delta s^{-1}\log(s)^{1/2}/3}\right]\right\}, \end{aligned}$$

which vanishes for any  $\delta$  large enough, as  $s^{-1}\log(s)^{1/2} \rightarrow 0$ .  $\square$

*Proof of Lemma 4(b).* For a sequence  $r_n \rightarrow \infty$  to be given later, define

$$H_i = s^{-1}(Kt)(X_{h,i})g(X_i)(Y_i\mathbb{1}\{Y_i \leq r_n\} - \mathbb{E}[Y_i\mathbb{1}\{Y_i \leq r_n\} | X_i])$$

and

$$T_i = s^{-1}(Kt)(X_{h,i})g(X_i)(Y_i\mathbb{1}\{Y_i > r_n\} - \mathbb{E}[Y_i\mathbb{1}\{Y_i > r_n\} | X_i]).$$

By the conditions on  $g(\cdot)$  and  $t(\cdot)$  and the kernel function,

$$|H_i| < C_1s^{-1}r_n$$

and

$$\mathbb{V}[H_i] = s^{-2}\mathbb{V}[(Kt)(X_{h,i})g(X_i)Y_i\mathbb{1}\{Y_i \leq r_n\}] \leq s^{-2}\mathbb{E}[(Kt)(X_{h,i})^2g(X_i)^2Y_i^2\mathbb{1}\{Y_i \leq r_n\}]$$

$$\begin{aligned}
&\leq s^{-2} \mathbb{E} [(Kt)(X_{h,i})^2 g(X_i)^2 Y_i^2] \\
&= s^{-2} \int (Kt)(X_{h,i})^2 g(X_i)^2 v(X_i) f(X_i) dX_i \\
&= s^{-2} h \int (Kt)(u)^2 (g v f)(x - uh) du \\
&\leq C_2/n.
\end{aligned}$$

Therefore, by Bernstein's inequality

$$\begin{aligned}
s^2 \mathbb{P} \left[ \left| \sum_{i=1}^n H_i \right| > \delta \log(s)^{1/2} \right] &\leq 2s^2 \exp \left\{ -\frac{\delta^2 \log(s)/2}{C_2 + C_1 s^{-1} r_n \delta \log(s)^{1/2}/3} \right\} \\
&\leq 2 \exp\{2 \log(s)\} \exp \left\{ -\frac{\delta^2 \log(s)/2}{C_2 + C_1 s^{-1} r_n \delta \log(s)^{1/2}/3} \right\} \\
&\leq 2 \exp \left\{ \log(s) \left[ 2 - \frac{\delta^2/2}{C_2 + C_1 s^{-1} r_n \delta \log(s)^{1/2}/3} \right] \right\},
\end{aligned}$$

which vanishes for  $\delta$  large enough as long as  $s^{-1} r_n \log(s)^{1/2}$  does not diverge.

Next, by Markov's inequality and the moment condition on  $Y$  of Assumption [S.II.3.1](#)

$$\begin{aligned}
s^2 \mathbb{P} \left[ \left| \sum_{i=1}^n T_i \right| > \delta \log(s)^{1/2} \right] &\leq s^2 \frac{1}{\delta^2 \log(s)} \mathbb{E} \left[ \left| \sum_{i=1}^n T_i \right|^2 \right] \\
&\leq s^2 \frac{1}{\delta^2 \log(s)} n \mathbb{E} [T_i^2] \\
&\leq s^2 \frac{1}{\delta^2 \log(s)} n \mathbb{V} [s^{-1} (Kt)(X_{h,i}) g(X_i) Y_i \mathbb{1}\{Y_i > r_n\}] \\
&\leq s^2 \frac{1}{\delta^2 \log(s)} n s^{-2} \mathbb{E} [(Kt)(X_{h,i})^2 g(X_i)^2 Y_i^2 \mathbb{1}\{Y_i > r_n\}] \\
&\leq s^2 \frac{1}{\delta^2 \log(s)} n s^{-2} \mathbb{E} [(Kt)(X_{h,i})^2 g(X_i)^2 |Y_i|^{2+\xi} r_n^{-\eta}] \\
&\leq s^2 \frac{1}{\delta^2 \log(s)} n s^{-2} (C h r_n^{-\xi}) \\
&\leq \frac{C}{\delta^2 \log(s)} \frac{s^2}{r_n^\xi},
\end{aligned}$$

which vanishes if  $s^2 \log(s)^{-1} r_n^{-\xi} \rightarrow 0$ .

It thus remains to choose  $r_n$  such that  $s^{-1} r_n \log(s)^{1/2}$  does not diverge and  $s^2 \log(s)^{-1} r_n^{-\xi} \rightarrow 0$ . This can be accomplished by setting  $r_n = s^\gamma$  for any  $2/\xi \leq \gamma < 1$ , which is possible as  $\xi > 2$ .  $\square$



*Proof of Lemma 4(c).* By Markov's inequality

$$\begin{aligned}
& \frac{1}{h^{p+1}} \mathbb{P} \left[ \left| s^{-2} \sum_{i=1}^n (Kt)(X_{h,i}) g(X_i) [m(X_i) - r_p(X_i - x)' \beta_p]^k \right| > \delta h^{(k-1)(p+1)} \log(s)^\gamma \right] \\
& \leq \frac{1}{h^{p+1}} \frac{1}{\delta h^{(k-1)(p+1)} \log(s)^\gamma} \mathbb{E} \left[ h^{-1} (Kt)(X_{h,i}) g(X_i) [m(X_i) - r_p(X_i - x)' \beta_p]^k \right] \\
& \leq \frac{1}{\delta h^{k(p+1)} \log(s)^\gamma} h^{k(p+1)} \mathbb{E} \left[ h^{-1} (Kt)(X_{h,i}) g(X_i) [h^{-p-1} (m(X_i) - r_p(X_i - x)' \beta_p)]^k \right] \\
& = O(\log(s)^{-\gamma}) = o(1).
\end{aligned}$$

This relies on the following calculation, which uses the conditions placed on  $m(\cdot)$ :

$$\begin{aligned}
& \mathbb{E} \left[ h^{-1} ((Kt)(X_{h,i}) g(X_i) \varepsilon_i) [m(X_i) - r_p(X_i - x)' \beta_p]^k \right] \\
& = h^{-1} \int (gfv)(X_i) (Kt)(X_{h,i}) [m(X_i) - r_p(X_i - x)' \beta_p]^k dX_i \\
& = h^{-1} \int (gfv)(X_i) (Kt)(X_{h,i}) \left( \frac{m^{(p+1)}(\bar{x})}{(p+1)!} (X_i - x)^{p+1} \right)^k dX_i \\
& = h^{k(p+1)} h^{-1} \int (gfv)(X_i) (Kt)(X_{h,i}) \left( \frac{m^{(p+1)}(\bar{x})}{(p+1)!} X_{h,i}^{p+1} \right)^k dX_i \\
& = Ch^{k(p+1)} h^{-1} \int (gfv)(X_i) (Kt)(X_{h,i}) X_{h,i}^{k(p+1)} dX_i \\
& = Ch^{k(p+1)} \int (gfv)(x - uh) (Kt)(u) u^{k(p+1)} du \\
& \asymp h^{k(p+1)}.
\end{aligned}$$

□

*Proof of Lemma 4(d).* By Markov's inequality, since  $\varepsilon_i$  is conditionally mean zero, we have

$$\begin{aligned}
& s^2 \mathbb{P} \left[ \left| s^{-2} \sum_{i=1}^n (Kt)(X_{h,i}) g(X_i) \varepsilon_i [m(X_i) - r_p(X_i - x)' \beta_p] \right| > \delta h^{p+1} \log(s)^\gamma \right] \\
& \leq s^2 \frac{1}{\delta h^{2(p+1)} \log(s)^{2\gamma}} \frac{1}{s^2} \mathbb{E} \left[ h^{-1} ((Kt)(X_{h,i}) g(X_i) \varepsilon_i)^2 [m(X_i) - r_p(X_i - x)' \beta_p]^2 \right] \\
& \leq \frac{s^2 h^{2(p+1)}}{\delta s^2 h^{2(p+1)} \log(s)^\gamma} \mathbb{E} \left[ h^{-1} ((Kt)(X_{h,i}) g(X_i) \varepsilon_i)^2 [h^{-p-1} (m(X_i) - r_p(X_i - x)' \beta_p)]^2 \right] \\
& \asymp \log(s)^{-2\gamma} \rightarrow 0,
\end{aligned}$$

where we rely on the same argument as above to compute the bias rate. □

*Proof of Lemma 4(e).* Follows from identical steps to 4(d). □

To illustrate how the above Lemma is used for the objects under study, we present the following collection of results. This is not meant to be an exhaustive list of all such results needed to prove all parts of Theorem 5, but any and all omitted terms follow by identical reasoning.

**Lemma 5.** *Let the conditions of Theorem 5 hold.*

(a) *For some  $\delta > 0$ ,  $r_*^{-1}\mathbb{P}[|\Gamma_p - \tilde{\Gamma}_p| > s^{-1}\log(s)^{1/2}] \rightarrow 0$ . Consequently, there exists a constant  $C_\Gamma < \infty$  such that  $\mathbb{P}[\Gamma_p^{-1} > 2C_\Gamma] = o(s^{-2})$  and so the prior rate result holds for  $|\Gamma_p^{-1} - \tilde{\Gamma}_p^{-1}|$  as well. Finally, these same results hold for  $\Gamma_q$  as well.*

(b) *For some  $\delta > 0$ ,  $r_*^{-1}\mathbb{P}[|\Lambda_{p,1} - \tilde{\Lambda}_{p,1}| > s^{-1}\log(s)^{1/2}] \rightarrow 0$ .*

(c) *For some  $\delta > 0$ ,*

$$s^2\mathbb{P}\left[\left|s^{-1}\sum_{i=1}^n\{(Kr_p)(X_{h,i})\varepsilon_i\}\right| > \delta\log(s)^{1/2}\right] \rightarrow 0.$$

(d) *For any  $\delta > 0$  and  $\gamma > 0$ ,*

$$\frac{1}{h^{p+1}}\mathbb{P}\left[\left|s^{-2}\sum_{i=1}^n\{(Kr_p)(X_{h,i})[m(X_i) - r_p(X_i - x)'\beta_p]\}\right| > \delta\log(s)^\gamma\right] \rightarrow 0.$$

(e) *There is some constant  $C_\Psi$  such that  $\mathbb{P}[\tilde{\Psi}_p > 2C_\Psi] = o(s^{-2})$ .*

*Proof of Lemma 5(a).* A typical element of  $\Gamma_p - \tilde{\Gamma}_p$  is, for some integer  $k \leq 2p$ ,

$$\frac{1}{nh}\sum_{i=1}^n\{K(X_{h,i})\mathcal{X}_{h,i}^k - \mathbb{E}[K(X_{h,i})\mathcal{X}_{h,i}^k]\}.$$

Therefore, the result follows by applying Lemma 4(a) to each element. Next, note that under the maintained assumptions

$$\tilde{\Gamma}_p = \mathbb{E}[h^{-1}(Kr_p r_p')(X_{h,i})] = h^{-1}\int(Kr_p r_p')(X_{h,i})f(X_i)dX_i = \int(Kr_p r_p')(u)f(x - uh)du$$

is bounded away from zero and infinity for  $n$  large enough. Therefore, there is a  $C_\Gamma < \infty$  such that  $|\tilde{\Gamma}_p^{-1}| < C_\Gamma$  and then

$$\begin{aligned}\mathbb{P}[\Gamma_p^{-1} > 2C_\Gamma] &= \mathbb{P}\left[\left(\Gamma_p^{-1} - \tilde{\Gamma}_p^{-1}\right) + \tilde{\Gamma}_p^{-1} > 2C_\Gamma\right] \\ &\leq \mathbb{P}\left[\Gamma_p^{-1} - \tilde{\Gamma}_p^{-1} > s^{-1}\log(s)^{1/2}\right] + \mathbb{P}\left[\tilde{\Gamma}_p^{-1} > 2C_\Gamma - s^{-1}\log(s)^{1/2}\right]\end{aligned}$$

$$= o(s^{-2}).$$

The third result follows from these two and the identity  $\Gamma_p^{-1} - \tilde{\Gamma}_p^{-1} = \tilde{\Gamma}_p^{-1}(\tilde{\Gamma}_p - \Gamma_p)\Gamma_p^{-1}$ .

Finally, for  $\Gamma_q$ , the identical steps apply with  $L$ ,  $q$ , and  $b$  in place of  $K$ ,  $p$ , and  $h$ .  $\square$

*Proof of Lemma 5(b).* Follows from identical steps to the previous result.  $\square$

*Proof of Lemma 5(c).* Follows from identical steps, but using Lemma 4(b) in place of Lemma 4(a).  $\square$

*Proof of Lemma 5(d).* Follows from identical steps, but using Lemma 4(c) in place of Lemma 4(a).  $\square$

*Proof of Lemma 5(e).* A typical element of  $\check{\Psi}_p$  is

$$\frac{1}{nh} \sum_{i=1}^n (K^2 r_p r'_p)(X_{h,i}) v(X_i),$$

and hence under the maintained assumptions the result follows just as the comparable result on  $\Gamma_p$ .  $\square$

We next state, without proof, the following fact about the rates appearing in all these Lemmas, which follows from elementary inequalities.

**Lemma 6.** *If  $r_1 = O(r'_1)$  and  $r_2 = O(r'_2)$ , for sequences of positive numbers  $r_1$ ,  $r'_1$ ,  $r_2$ , and  $r'_2$  and if a sequence of nonnegative random variables obeys  $(r_1)^{-1} \mathbb{P}[U_n > r_2] \rightarrow 0$  it also holds that  $(r'_1)^{-1} \mathbb{P}[U_n > r'_2] \rightarrow 0$ .*

*In particular, since  $r_* = \max\{s^{-2}, \eta^2, s^{-1}\eta\}$  is defined as the slowest vanishing of the rates, then  $r_1^{-1} \mathbb{P}[|U'| > r_n] = o(1)$  implies  $r_*^{-1} \mathbb{P}[|U'| > r_n] = o(1)$ , for  $r_1$  equal to any of  $s^{-2}$ ,  $\eta^2$ , or  $s^{-1}\eta$ . Similarly,  $r_n$  may be chosen as any sequence that obeys  $r_n = o(r_*)$ . Thus, for different pieces of  $U$  defined in Eqn. (29), we may make different choices for these two sequences, as convenient.*

The next Lemma proves Eqn. (29), a crucial step in the proof of Theorem 5(a). Because this result only involves undersmoothing, we will omit the subscript  $p$  as above.

**Lemma 7.** *Let the conditions of Theorem 5(a) hold. Then Eqn. (29) holds, namely, for some  $r_n = o(r_*)$*

$$\frac{1}{r_*} \mathbb{P}[|U| > r_n] \rightarrow 0.$$

*Proof.* Recall the definition:

$$\begin{aligned}
U = & \left\{ -\frac{1}{2\tilde{\sigma}^3} e'_0 \Gamma^{-1} (A_{1,6} + A_{1,7} + A_{1,8}) \Gamma^{-1} e_0 + \frac{3}{8\tilde{\sigma}^5} \left[ e'_0 \Gamma^{-1} \left( \sum_{k=2}^8 A_{1,k} \right) \Gamma^{-1} e_0 \right]^2 - \frac{5}{16} \frac{(\hat{\sigma}^2 - \tilde{\sigma}^2)^3}{\tilde{\sigma}^7} \right\} \\
& \times \left\{ s e'_0 \Gamma^{-1} R'W(Y - M)/n + s e'_0 \Gamma^{-1} R'W(M - \check{R}\beta)/n \right\} \\
& - \left\{ -\frac{1}{2\tilde{\sigma}^3} e'_0 \tilde{\Gamma}^{-1} \left( \tilde{A}_{1,6} + \tilde{A}_{1,7} + \tilde{A}_{1,8} \right) \tilde{\Gamma}^{-1} e_0 \right\} \eta.
\end{aligned}$$

To fully prove the claim of the lemma, we must fully expand  $U$  and bound each piece. First, we present complete details on two terms. The remainder are entirely analogous, as discussed below. Consider the pieces involving  $A_{1,6}$ , namely:

$$e'_0 \Gamma^{-1} A_{1,6} \Gamma^{-1} e_0 \left\{ s e'_0 \Gamma^{-1} R'W(Y - M)/n + s e'_0 \Gamma^{-1} R'W(M - \check{R}\beta)/n \right\} - e'_0 \tilde{\Gamma}^{-1} \tilde{A}_{1,6} \tilde{\Gamma}^{-1} e_0 \eta.$$

The first of these is

$$\begin{aligned}
e'_0 \Gamma^{-1} A_{1,6} \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} R'W(Y - M)/n &= e'_0 \Gamma^{-1} \left( A_{1,6} - \tilde{A}_{1,6} \right) \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} R'W(Y - M)/n \\
&+ e'_0 \left( \Gamma^{-1} - \tilde{\Gamma}^{-1} \right) \tilde{A}_{1,6} \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} R'W(Y - M)/n \\
&+ e'_0 \tilde{\Gamma}^{-1} \tilde{A}_{1,6} \left( \Gamma^{-1} - \tilde{\Gamma}^{-1} \right) e_0 s e'_0 \Gamma^{-1} R'W(Y - M)/n \\
&+ e'_0 \tilde{\Gamma}^{-1} \tilde{A}_{1,6} \tilde{\Gamma}^{-1} e_0 s e'_0 \left( \Gamma^{-1} - \tilde{\Gamma}^{-1} \right) R'W(Y - M)/n \\
&+ e'_0 \tilde{\Gamma}^{-1} \tilde{A}_{1,6} \tilde{\Gamma}^{-1} e_0 s e'_0 \tilde{\Gamma}^{-1} R'W(Y - M)/n. \\
&=: U_{1,1} + U_{1,2} + U_{1,3} + U_{1,4} + U_{1,5}
\end{aligned}$$

We now bound each remainder in turn. First, for  $r_n = h^{p+1} \log(s)^{-1/2}$ , we have

$$\begin{aligned}
s^2 \mathbb{P}[|U_{1,1}| > r_n] &= s^2 \mathbb{P} \left[ \left| e'_0 \Gamma^{-1} \left( A_{1,6} - \tilde{A}_{1,6} \right) \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} R'W(Y - M)/n \right| > r_n \right] \\
&\leq s^2 \mathbb{P} \left[ 8C_\Gamma^3 \left| A_{1,6} - \tilde{A}_{1,6} \right| > \log(s)^{-1/2} r_n \right] \\
&\quad + s^2 \mathbb{P} \left[ \left| s^{-1} \sum_{i=1}^n \{ (K r_p)(X_{h,i}) \varepsilon_i \} \right| > \log(s)^{1/2} \right] + s^2 3 \mathbb{P}[\Gamma_p^{-1} > 2C_\Gamma] \\
&= s^2 \mathbb{P} \left[ 8C_\Gamma^3 \left| A_{1,6} - \tilde{A}_{1,6} \right| > h^{2(p+1)} \log(s)^\gamma \frac{r_n}{h^{2(p+1)} \log(s)^{1/2+\gamma}} \right] + o(1) \\
&= o(1),
\end{aligned}$$

because  $h^{-2(p+1)} r_n \log(s)^{-1/2-\gamma} = h^{-(p+1)} \log(s)^{-1-\gamma} \rightarrow \infty$ .

Next, since  $\tilde{A}_{1,6} \asymp h^{2(p+1)}$ , for  $r_n = h^{p+1} \log(s)^{-1/2}$ .

$$\begin{aligned}
s^2 \mathbb{P}[|U_{1,2}| > r_n] &= s^2 \mathbb{P}\left[\left|e'_0 \left(\Gamma^{-1} - \tilde{\Gamma}^{-1}\right) \tilde{A}_{1,6} \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} R' W(Y - M)/n\right| > r_n\right] \\
&\leq s^2 \mathbb{P}\left[4C_\Gamma^2 \left|\tilde{A}_{1,6}\right| \left|s^{-1} \sum_{i=1}^n \{(Kr_p)(X_{h,i})\varepsilon_i\}\right| > s \log(s)^{-1/2} r_n\right] \\
&\quad + s^2 \mathbb{P}\left[\left|\Gamma^{-1} - \tilde{\Gamma}^{-1}\right| > s^{-1} \log(s)^{1/2}\right] + s^2 2\mathbb{P}\left[\Gamma_p^{-1} > 2C_\Gamma\right] \\
&= s^2 \mathbb{P}\left[4C_\Gamma^2 \left|s^{-1} \sum_{i=1}^n \{(Kr_p)(X_{h,i})\varepsilon_i\}\right| > \log(s)^{1/2} \frac{sr_n}{h^{2(p+1)} \log(s)}\right] + o(1) \\
&= o(1),
\end{aligned}$$

because  $sr_n h^{-2(p+1)} \log(s)^{-1} = sh^{-(p+1)} \log(s)^{-3/2} \rightarrow \infty$ . Terms  $U_{1,3}$  and  $U_{1,4}$  are nearly identically treated.

Let  $r_n = h^{p+1} \log(s)^{-1/2}$ . Then since  $\tilde{A}_{1,6} \asymp h^{2(p+1)}$ ,

$$\begin{aligned}
s^2 \mathbb{P}[|U_{1,5}| > r_n] &= s^2 \mathbb{P}\left[\left|e'_0 \tilde{\Gamma}^{-1} \tilde{A}_{1,6} \tilde{\Gamma}^{-1} e_0 s e'_0 \tilde{\Gamma}^{-1} R' W(Y - M)/n\right| > r_n\right] \\
&\leq s^2 \mathbb{P}\left[C_\Gamma^3 \left|\tilde{A}_{1,6}\right| \left|s^{-1} \sum_{i=1}^n \{(Kr_p)(X_{h,i})\varepsilon_i\}\right| > r_n\right] \\
&\leq s^2 \mathbb{P}\left[C_\Gamma^3 \left|s^{-1} \sum_{i=1}^n \{(Kt)(X_{h,i})g(X_i)\varepsilon_i\}\right| > \log(s)^{1/2} \frac{\log(s)^{-1/2} r_n}{h^{2(p+1)}}\right] \\
&= o(1),
\end{aligned}$$

because  $h^{-2(p+1)} r_n \log(s)^{-1/2} = h^{-(p+1)} \log(s)^{-1} \rightarrow \infty$ .

Thus, since  $\tilde{\sigma}^{-1}$  is bounded away from zero, we find that

$$s^2 \mathbb{P}\left[\left|\frac{1}{2\tilde{\sigma}^3} e'_0 \Gamma^{-1} A_{1,6} \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} R' W(Y - M)/n\right| > r_n\right] \rightarrow 0.$$

Turning our attention to the second term, we have

$$\begin{aligned}
&e'_0 \Gamma^{-1} A_{1,6} \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} R' W(M - \check{R}\beta)/n - e'_0 \tilde{\Gamma}^{-1} \tilde{A}_{1,6} \tilde{\Gamma}^{-1} e_0 \eta \\
&= e'_0 \Gamma^{-1} \left(A_{1,6} - \tilde{A}_{1,6}\right) \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} R' W(M - \check{R}\beta)/n \\
&\quad + e'_0 \Gamma^{-1} \tilde{A}_{1,6} \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} \left(R' W(M - \check{R}\beta)/n - \mathbb{E}[R' W(M - \check{R}\beta)/n]\right) \\
&\quad + e'_0 \left(\Gamma^{-1} - \tilde{\Gamma}^{-1}\right) \tilde{A}_{1,6} \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} \mathbb{E}[R' W(M - \check{R}\beta)/n] \\
&\quad + e'_0 \tilde{\Gamma}^{-1} \tilde{A}_{1,6} \left(\Gamma^{-1} - \tilde{\Gamma}^{-1}\right) e_0 s e'_0 \Gamma^{-1} \mathbb{E}[R' W(M - \check{R}\beta)/n]
\end{aligned}$$

$$\begin{aligned}
& + e'_0 \tilde{\Gamma}^{-1} \tilde{A}_{1,6} \tilde{\Gamma}^{-1} e_0 s e'_0 \left( \Gamma^{-1} - \tilde{\Gamma}^{-1} \right) \mathbb{E} \left[ R'W(M - \check{R}\beta)/n \right] \\
& =: U_{2,1} + U_{2,2} + U_{2,3} + U_{2,4} + U_{2,5}.
\end{aligned}$$

For  $r_n = h^{p+1} \log(s)^{-1}$ , we have

$$\begin{aligned}
r_*^{-1} \mathbb{P} [|U_{2,1}| > r_n] &= r_*^{-1} \mathbb{P} \left[ e'_0 \Gamma^{-1} \left( A_{1,6} - \tilde{A}_{1,6} \right) \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} R'W(M - \check{R}\beta)/n > r_n \right] \\
&\leq r_*^{-1} \mathbb{P} \left[ 8C_\Gamma^3 s \left| A_{1,6} - \tilde{A}_{1,6} \right| > sh^{2(p+1)} \log(s)^\gamma \frac{r_n}{sh^{2(p+1)} \log(s)^{2\gamma}} \right] \\
&\quad + r_*^{-1} \mathbb{P} \left[ \left| \frac{1}{nh} \sum_{i=1}^n \{ (Kr_p)(X_{h,i}) [m(X_i) - r_p(X_i - x)' \beta_p] \} \right| > \log(s)^\gamma \right] \\
&\quad + r_*^{-1} 3 \mathbb{P} [\Gamma_p^{-1} > 2C_\Gamma] \\
&\leq s^2 \mathbb{P} \left[ 8C_\Gamma^3 s \left| A_{1,6} - \tilde{A}_{1,6} \right| > sh^{2(p+1)} \log(s)^\gamma \frac{r_n}{sh^{2(p+1)} \log(s)^{2\gamma}} \right] \\
&\quad + h^{-(p+1)} \mathbb{P} \left[ \left| \frac{1}{nh} \sum_{i=1}^n \{ (Kr_p)(X_{h,i}) [m(X_i) - r_p(X_i - x)' \beta_p] \} \right| > \log(s)^\gamma \right] \\
&\quad + s^2 3 \mathbb{P} [\Gamma_p^{-1} > 2C_\Gamma] \\
&= o(1),
\end{aligned}$$

because  $sh^{2(p+1)} r_n^{-1} \log(s)^{2\gamma} = sh^{p+1} \log(s)^{1+2\gamma} \rightarrow 0$  by the conditions on  $\eta$  placed in the theorem.

Next, with  $r_n = h^{p+1} \log(s)^{-1}$  and using  $\tilde{A}_{1,6} \asymp h^{2(p+1)}$ , we have

$$\begin{aligned}
r_*^{-1} \mathbb{P} [|U_{2,2}| > r_n] &= r_*^{-1} \mathbb{P} \left[ \left| e'_0 \Gamma^{-1} \tilde{A}_{1,6} \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} (R'W(M - \check{R}\beta)/n - \mathbb{E} [R'W(M - \check{R}\beta)/n]) \right| > r_n \right] \\
&\leq r_*^{-1} \mathbb{P} \left[ 8C_\Gamma^3 \left| \tilde{A}_{1,6} \right| \left| s^{-1} \sum_{i=1}^n \left\{ (Kr_p)(X_{h,i}) [m(X_i) - r_p(X_i - x)' \beta_p] \right. \right. \right. \\
&\quad \left. \left. \left. - \mathbb{E} [(Kr_p)(X_{h,i}) [m(X_i) - r_p(X_i - x)' \beta_p]] \right\} \right| > r_n \right] \\
&\quad + r_*^{-1} 3 \mathbb{P} [\Gamma_p^{-1} > 2C_\Gamma] \\
&\leq s^2 \mathbb{P} \left[ 8C_\Gamma^3 \left| s^{-2} \sum_{i=1}^n \left\{ (Kr_p)(X_{h,i}) [m(X_i) - r_p(X_i - x)' \beta_p] \right. \right. \right. \\
&\quad \left. \left. \left. - \mathbb{E} [(Kr_p)(X_{h,i}) [m(X_i) - r_p(X_i - x)' \beta_p]] \right\} \right| > h^{p+1} \log(s)^\gamma \frac{r_n}{h^{3(p+1)} \log(s)^\gamma} \right] \\
&\quad + s^2 3 \mathbb{P} [\Gamma_p^{-1} > 2C_\Gamma] \\
&= o(1),
\end{aligned}$$

because  $r_n h^{-3(p+1)} \log(s)^{-\gamma} = h^{-2(p+1)} \log(s)^{-1-\gamma} \rightarrow \infty$ .

Third, as  $\tilde{A}_{1,6} \asymp h^{2(p+1)}$  and  $\mathbb{E} [R'W(M - \tilde{R}\beta)/n] \asymp h^{p+1}$ , if we choose  $r_n = h^{p+1} \log(s)^{-1}$ ,

$$\begin{aligned} r_*^{-1} \mathbb{P} [|U_{2,3}| > r_n] &\leq r_*^{-1} \mathbb{P} \left[ 4C_\Gamma^2 s \left| \Gamma^{-1} - \tilde{\Gamma}^{-1} \right| > s^{-1} \log(s)^{1/2} \frac{s r_n}{h^{3(p+1)} \log(s)^{1/2}} \right] \\ &\quad + r_*^{-1} 2 \mathbb{P} [\Gamma_p^{-1} > 2C_\Gamma] \\ &\leq s^2 \mathbb{P} \left[ 4C_\Gamma^2 \left| \Gamma^{-1} - \tilde{\Gamma}^{-1} \right| > s^{-1} \log(s)^{1/2} \frac{r_n}{h^{3(p+1)} \log(s)^{1/2}} \right] \\ &\quad + s^2 2 \mathbb{P} [\Gamma_p^{-1} > 2C_\Gamma] \\ &= o(1), \end{aligned}$$

because  $r_n h^{-3(p+1)} \log(s)^{-1/2} = h^{-2(p+1)} \log(s)^{-1-1/2} \rightarrow \infty$ . The terms  $U_{2,3}$  and  $U_{2,5}$  are handled identically.

Thus, since  $\tilde{\sigma}^{-1}$  is bounded away from zero, we find that

$$s^2 \mathbb{P} \left[ \left| \frac{1}{2\tilde{\sigma}^3} e'_0 \Gamma^{-1} A_{1,6} \Gamma^{-1} e_0 s e'_0 \Gamma^{-1} R'W(M - \tilde{R}\beta)/n - e'_0 \tilde{\Gamma}^{-1} \tilde{A}_{1,6} \tilde{\Gamma}^{-1} e_0 \eta \right| > r_n \right] \rightarrow 0.$$

The same type of arguments, though notationally more challenging, will show that the remainder of  $U$  obeys the same bounds. Note that the rest of the terms are even higher order, involving either  $A_{1,7}$  and  $A_{1,8}$ , or the square or cube of the other errors. It is for this reason that only the “leading” three terms need be centered, that is, why only

$$- \left\{ -\frac{1}{2\tilde{\sigma}^3} e'_0 \tilde{\Gamma}^{-1} \left( \tilde{A}_{1,6} + \tilde{A}_{1,7} + \tilde{A}_{1,8} \right) \tilde{\Gamma}^{-1} e_0 \right\} \eta$$

appears in  $\tilde{z}$ . □

## S.II.6.4 Computing the Terms of the Expansion

Identifying the terms of the expansion is a matter of straightforward, if tedious, calculation. The first four cumulants of the Studentized statistics must be calculated (due to [James and Mayne \(1962\)](#)), which are functions of the first four moments. In what follows, we give a short summary. Note well that we always discard higher-order terms for brevity, and to save notation we will write  $\overset{o}{=}$  to stand in for “equal up to  $o((nh)^{-1} + (nh)^{-1/2} \eta + \eta^2)$ ”, and including  $o(\rho^{1+2(p+1)})$  for  $T_{bc}$ .

The computations will be aided by putting all three estimators into a common structure. In close parallel to the density case, let us define  $\hat{m}_1 := \hat{m}$  and  $\hat{m}_2 = \hat{m} - \hat{m}_m$ ,  $\sigma_1^2 := \sigma_{us}^2$ , and  $\sigma_2^2 := \sigma_{rbc}^2$ , so that subscripts 1 and 2 generically stand in for undersmoothing and bias

correction, respectively. With this in mind, we write

$$T_{\text{us}} = T_{1,1}, \quad T_{\text{bc}} = T_{2,1}, \quad \text{and} \quad T_{\text{rbc}} = T_{2,2},$$

again paralleling the density case, so that the first subscript refers to the numerator and the second to the denominator. In the same vein, with some abuse of notation, we will also use<sup>3</sup>  $r_1(u) = r_p(u)$ ,  $r_2(u) = r_q(u)$ ,  $K_1(u) = K(u)$ ,  $K_2(u) = L(u)$ ,  $h_1 = h$ , and  $h_2 = b$ , as well as

$$\begin{aligned} \ell_1^0(X_i) &\equiv \ell_{\text{us}}^0(X_i), \\ \ell_1^1(X_i, X_j) &\equiv \ell_{\text{us}}^1(X_i, X_j), \\ \ell_2^0(X_i) &\equiv \ell_{\text{bc}}^0(X_i), \\ \ell_2^1(X_i, X_j) &\equiv \ell_{\text{bc}}^1(X_i, X_j). \end{aligned}$$

For the purpose of computing the expansion terms (i.e. moments of the two sides agree up to the requisite order), recalling the Taylor series expansion above, we will use

$$T_{v,w} \approx \left\{ 1 - \frac{1}{2\tilde{\sigma}_w^2} (W_{w,1} + V_{w,1} + V_{w,2}) + \frac{3}{8\tilde{\sigma}_w^4} (W_{w,1} + V_{w,1} + V_{w,2})^2 \right\} \tilde{\sigma}_w^{-1} \{E_{v,1} + E_{v,2} + E_{v,3} + B_{v,1}\},$$

where we define, for  $v \in \{1, 2\}$ ,

$$\begin{aligned} E_{v,1} &= s \frac{1}{nh} \sum_{i=1}^n \ell_v^0(X_i) \varepsilon_i \\ E_{v,2} &= s \frac{1}{(nh)^2} \sum_{i=1}^n \sum_{j=1}^n \ell_v^1(X_i, X_j) \varepsilon_i, \\ E_{v,3} &=: s \frac{1}{(nh)^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \ell_v^2(X_i, X_j, X_k) \varepsilon_i, \end{aligned}$$

where the final line defines  $\ell_{\text{us}}^2(X_i, X_j, X_k)$  in the obvious way following  $\ell_{\text{us}}^1$ . To concretize the notation, for undersmoothing we are defining

$$\begin{aligned} E_{1,1} &= s e_0' \tilde{\Gamma}_p^{-1} R_p' W_p (Y - M) / n, \\ E_{1,2} &= s e_0' \tilde{\Gamma}_p^{-1} (\tilde{\Gamma}_p - \Gamma_p) \tilde{\Gamma}_p^{-1} R_p' W_p (Y - M) / n, \\ E_{1,3} &= s e_0' \tilde{\Gamma}_p^{-1} (\tilde{\Gamma}_p - \Gamma_p) \tilde{\Gamma}_p^{-1} (\tilde{\Gamma}_p - \Gamma_p) \tilde{\Gamma}_p^{-1} R_p' W_p (Y - M) / n. \end{aligned}$$

---

<sup>3</sup>Throughout Section [S.II](#), we use only generic polynomial orders  $p$  and  $q$ , and so this notation will not conflict with the local linear or local quadratic fits, which would also be denoted  $r_1(u)$  and  $r_2(u)$ , respectively.



In a similar way,

$$\begin{aligned}
W_{v,1} &= \frac{1}{nh} \sum_{i=1}^n \left\{ \ell_v^0(X_i)^2 (\varepsilon_i^2 - v(X_i)) \right\} - 2 \frac{1}{n^2 h^2} \sum_{i=1}^n \sum_{j=1}^n \left\{ \ell_v^0(X_i)^2 r_v(X_{h_v,i})' \tilde{\Gamma}_v^{-1}(K_v r_v)(X_{h_v,i}) \varepsilon_i \varepsilon_j \right\} \\
&\quad + \frac{1}{n^3 h^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \left\{ \ell_v^0(X_i)^2 r_v(X_{h_v,i})' \tilde{\Gamma}_v^{-1}(K_v r_v)(X_{h_v,i}) \varepsilon_j \varepsilon_k \right\}, \\
V_{v,1} &= \frac{1}{nh} \sum_{i=1}^n \left\{ \ell_v^0(X_i)^2 v(X_i)^2 - \mathbb{E}[\ell_v^0(X_i)^2 v(X_i)^2] \right\} + 2 \frac{1}{n^2 h^2} \sum_{i=1}^n \sum_{j=1}^n \ell_v^2(X_i, X_j) \ell_v^0(X_i) v(X_i), \\
V_{v,2} &= \frac{1}{n^3 h^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \ell_v^1(X_i, X_j) \ell_v^1(X_i, X_k) v(X_i) + 2 \frac{1}{n^3 h^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \ell_v^2(X_i, X_j, X_k) \ell_v^0(X_i) v(X_i),
\end{aligned}$$

and specifically for undersmoothing and bias correction, let

$$B_{1,1} = s \frac{1}{nh} \sum_{i=1}^n \ell_1^0(X_i) [m(X_i) - r_p(X_i - x)' \beta_p]$$

and

$$\begin{aligned}
B_{2,1} &= s \frac{1}{nh} \sum_{i=1}^n \left\{ h^{-1} \ell_{\text{us}}^0(X_i) [m(X_i) - r_{p+1}(X_i - x)' \beta_{p+1}] \right. \\
&\quad \left. - h^{-1} (\ell_{\text{bc}}^0(X_i) - \ell_{\text{us}}^0(X_i)) [m(X_i) - r_q(X_i - x)' \beta_q] \right\}.
\end{aligned}$$

Note that  $\eta_{\text{us}} = \mathbb{E}[B_{1,1}]$  and  $\eta_{\text{bc}} = \mathbb{E}[B_{2,1}]$ .

Straightforward moment calculations yield

$$\mathbb{E}[T_{v,w}] \stackrel{o}{=} \tilde{\sigma}_w^{-1} \mathbb{E}[B_{v,1}] - \frac{1}{2\tilde{\sigma}_w^2} \mathbb{E}[W_{w,1} E_{v,1}],$$

$$\begin{aligned}
\mathbb{E}[T_{v,w}^2] &\stackrel{o}{=} \frac{1}{\tilde{\sigma}_w^2} \mathbb{E}[E_{v,1}^2 + E_{v,2}^2 + 2E_{v,1}E_{v,2} + 2E_{v,1}E_{v,3}] \\
&\quad - \frac{1}{\tilde{\sigma}_w^4} \mathbb{E}[W_{w,1}E_{v,1}^2 + V_{w,1}E_{v,1}^2 + V_{w,2}E_{v,1}^2 + 2V_{w,1}E_{v,1}E_{v,2}] \\
&\quad + \frac{1}{\tilde{\sigma}_w^6} \mathbb{E}[W_{w,1}^2E_{v,1}^2 + V_{w,1}^2E_{v,1}^2] + \frac{1}{\tilde{\sigma}_w^2} \mathbb{E}[B_{v,1}^2] - \frac{1}{\tilde{\sigma}_w^4} \mathbb{E}[W_{w,1}E_{v,1}B_{v,1}],
\end{aligned}$$

$$\mathbb{E}[T_{v,w}^3] \stackrel{o}{=} \frac{1}{\tilde{\sigma}_w^3} \mathbb{E}[E_{v,1}^3] - \frac{3}{2\tilde{\sigma}_w^5} \mathbb{E}[W_{w,1}E_{v,1}^3] + \frac{3}{\tilde{\sigma}_w^3} \mathbb{E}[E_{v,1}^2 B_{v,1}],$$

and

$$\begin{aligned}
\mathbb{E}[T_{v,w}^4] &\stackrel{o}{=} \frac{1}{\tilde{\sigma}_w^4} \mathbb{E} [E_{v,1}^4 + 4E_{v,1}^3 E_{v,2} + 4E_{v,1}^3 E_{v,3} + 6E_{v,1}^2 E_{v,3}^2] \\
&\quad - \frac{2}{\tilde{\sigma}_w^6} \mathbb{E} [W_{w,1} E_{v,1}^4 + V_{w,1} E_{v,1}^4 + 4V_{w,1} E_{v,1}^3 E_{v,2} + V_{w,2} E_{v,1}] \\
&\quad + \frac{3}{\tilde{\sigma}_w^8} \mathbb{E} [W_{w,1}^2 E_{v,1}^4 + V_{w,1}^2 E_{v,1}^4] \\
&\quad + \frac{4}{\tilde{\sigma}_w^4} \mathbb{E} [E_{v,1}^3 B_{v,1}] - \frac{8}{\tilde{\sigma}_w^6} \mathbb{E} [W_{w,1} E_{v,1}^3 B_{v,1}] + \frac{6}{\tilde{\sigma}_w^4} \mathbb{E} [E_{v,1}^2 B_{v,1}^2].
\end{aligned}$$

Computing each term in turn, we have

$$\begin{aligned}
\mathbb{E}[B_{v,1}] &= \eta_v, \\
\mathbb{E}[W_{w,1} E_{v,1}] &\stackrel{o}{=} s^{-1} \mathbb{E} [h^{-1} \ell_w^0(X_i)^2 \ell_v^0(X_i) \varepsilon_i^3], \\
\mathbb{E}[E_{v,1}^2] &\stackrel{o}{=} \tilde{\sigma}_v^2, \\
\mathbb{E}[E_{v,1} E_{v,2}] &\stackrel{o}{=} s^{-2} \mathbb{E} [h^{-1} \ell_v^1(X_i, X_i) \ell_v^0(X_i) \varepsilon_i^2], \\
\mathbb{E}[E_{v,2}^2] &\stackrel{o}{=} s^{-1} \mathbb{E} [h^{-2} \ell_v^1(X_i, X_j)^2 \varepsilon_i^2], \\
\mathbb{E}[E_{v,2} E_{v,3}] &\stackrel{o}{=} s^{-2} \mathbb{E} [h^{-2} \ell_v^2(X_i, X_j, X_j) \ell_v^0(X_i) \varepsilon_i^2], \\
\mathbb{E}[W_{w,1} E_{v,1}^2] &\stackrel{o}{=} s^{-2} \left\{ \mathbb{E} [h^{-1} \ell_w^0(X_i)^2 \ell_v^0(X_i)^2 (\varepsilon_i^4 - v(X_i)^2)] \right. \\
&\quad - 2\tilde{\sigma}_v^2 \mathbb{E} [h^{-1} \ell_w^0(X_i)^2 r_w(X_{h_w,i})' \tilde{\Gamma}_w^{-1} (K_w r_w)(X_{h_w,i}) \varepsilon_i^2] \\
&\quad - 4\mathbb{E} [h^{-1} \ell_w^0(X_i)^2 \ell_v^0(X_i)^2 r_w(X_{h_w,i})' \tilde{\Gamma}_w^{-1} \varepsilon_i^2] \mathbb{E} [h^{-1} (K_w r_w)(X_{h_w,i}) \ell_v^0(X_i) \varepsilon_i^2] \\
&\quad + \tilde{\sigma}_v^2 \mathbb{E} \left[ h^{-2} \ell_w^0(X_i)^2 \left( r_w(X_{h_w,i})' \tilde{\Gamma}_w^{-1} (K_w r_w)(X_{h_w,j}) \right)^2 \varepsilon_j^2 \right] \\
&\quad \left. + \mathbb{E} \left[ h^{-1} \ell_{\text{us}}^0(X_j)^2 \left( \mathbb{E} [h^{-1} r_p(X_{h,j})' \tilde{\Gamma}_p^{-1} (K r_p)(X_{h,i}) \ell_{\text{us}}^0(X_i) \varepsilon_i^2 | X_j] \right)^2 \right] \right\}, \\
\mathbb{E}[V_{w,1} E_{v,1}^2] &\stackrel{o}{=} s^{-2} \left\{ \mathbb{E} [h^{-1} (\ell_w^0(X_i)^2 v(X_i) - \mathbb{E}[\ell_w^0(X_i)^2 v(X_i)]) \ell_v^0(X_i)^2 \varepsilon_i^2] \right. \\
&\quad \left. + 2\tilde{\sigma}_v^2 \mathbb{E} [h^{-1} \ell_w^1(X_i, X_i) \ell_w^0(X_i) v(X_i)] \right\}, \\
\mathbb{E}[V_{w,1} E_{v,1} E_{v,2}] &\stackrel{o}{=} s^{-2} \left\{ \mathbb{E} [h^{-2} (\ell_w^0(X_j)^2 v(X_j) - \mathbb{E}[\ell_w^0(X_j)^2 v(X_j)]) \ell_v^1(X_i, X_j) \ell_v^0(X_i) \varepsilon_i^2] \right. \\
&\quad \left. + 2\mathbb{E} [h^{-3} \ell_w^1(X_i, X_j) \ell_v^1(X_k, X_j) \ell_w^0(X_i) \ell_v^0(X_k) v(X_i) \varepsilon_k^2] \right\}, \\
\mathbb{E}[V_{w,2} E_{v,1}^2] &\stackrel{o}{=} s^{-2} \left\{ \tilde{\sigma}_v^2 \mathbb{E} [h^{-2} (\ell_w^1(X_i, X_j)^2 + 2\ell_w^2(X_i, X_j, X_j)) v(X_i)] \right\}, \\
\mathbb{E}[W_{w,1}^2 E_{v,1}^2] &\stackrel{o}{=} s^{-2} \left\{ \tilde{\sigma}_v^2 \mathbb{E} [h^{-1} \ell_w^0(X_i)^4 (\varepsilon_i^4 - v(X_i)^2)] + 2\mathbb{E} [h^{-1} \ell_v^0(X_i) \ell_w^0(X_i)^2 \varepsilon_i^3]^2 \right\},
\end{aligned}$$

$$\begin{aligned}
\mathbb{E} [V_{w,1}^2 E_{v,1}^2] &\stackrel{o}{=} s^{-2} \tilde{\sigma}_v^2 \left\{ \mathbb{E} \left[ h^{-1} (\ell_w^0(X_i)^2 v(X_i) - \mathbb{E}[\ell_w^0(X_i)^2 v(X_i)])^2 \right] \right. \\
&\quad + 4\mathbb{E} \left[ h^{-2} (\ell_w^0(X_i)^2 v(X_i) - \mathbb{E}[\ell_w^0(X_i)^2 v(X_i)]) \ell_w^1(X_j, X_i) \ell_w^0(X_j) v(X_j) \right] \\
&\quad \left. + 4\mathbb{E} \left[ h^{-3} \ell_w^1(X_i, X_j) \ell_w^0(X_i) v(X_i) \ell_w^1(X_k, X_j) \ell_w^0(X_k) v(X_k) \right] \right\}, \\
\mathbb{E} [W_{w,1} E_{v,1} B_{v,1}] &\stackrel{o}{=} \mathbb{E} [W_{w,1} E_{v,1}] \mathbb{E} [B_{v,1}], \\
\mathbb{E} [E_{v,1}^3] &\stackrel{o}{=} s^{-1} \mathbb{E} [h^{-1} \ell_v^0(X_i)^3 \varepsilon_i^3], \\
\mathbb{E} [W_{w,1} E_{v,1}^3] &\stackrel{o}{=} \mathbb{E} [E_{v,1}^2] \mathbb{E} [W_{w,1} E_{v,1}], \\
\mathbb{E} [E_{v,1}^4] &\stackrel{o}{=} 3\tilde{\sigma}_v^4 + s^{-2} \mathbb{E} [h^{-1} \ell_v^0(X_i)^4 \varepsilon_i^3], \\
\mathbb{E} [E_{v,1}^3 E_{v,2}] &\stackrel{o}{=} s^{-2} 6\tilde{\sigma}_v^2 \mathbb{E} [h^{-1} \ell_v^1(X_i, X_i) \ell_v^0(X_i) \varepsilon_i^2], \\
\mathbb{E} [E_{v,1}^3 E_{v,3}] &\stackrel{o}{=} s^{-2} 3\tilde{\sigma}_v^2 \mathbb{E} [h^{-2} \ell_v^2(X_i, X_j, X_j) \ell_v^0(X_i) \varepsilon_i^2], \\
\mathbb{E} [E_{v,1}^2 E_{v,2}^2] &\stackrel{o}{=} s^{-2} \left\{ \tilde{\sigma}_v^2 \mathbb{E} [h^{-2} \ell_v^1(X_i, X_j)^2 \varepsilon_i^2] + 2\mathbb{E} [h^{-3} \ell_v^1(X_i, X_j) \ell_v^1(X_k, X_j) \ell_v^0(X_i) \ell_v^0(X_k) \varepsilon_i^2 \varepsilon_k^2] \right\}, \\
\mathbb{E} [W_{w,1} E_{v,1}^4] &\stackrel{o}{=} s^{-2} \left\{ \mathbb{E} [h^{-1} \ell_w^0(X_i)^2 \ell_v^0(X_i) \varepsilon_i^3] \mathbb{E} [h^{-1} \ell_v^0(X_i)^3 \varepsilon_i^3] + 6\mathbb{E} [E_{v,1}^2] \mathbb{E} [W_{w,1} E_{v,1}^2] \right\}, \\
\mathbb{E} [V_{w,1} E_{v,1}^4] &\stackrel{o}{=} s^{-2} \tilde{\sigma}_v^2 6 \left\{ \mathbb{E} [h^{-1} (\ell_w^0(X_i)^2 v(X_i) - \mathbb{E}[\ell_w^0(X_i)^2 v(X_i)]) \ell_v^0(X_i)^2 \varepsilon_i^2] \right. \\
&\quad \left. + 2\mathbb{E} [h^{-2} \ell_w^1(X_i, X_j) \ell_w^0(X_i) \ell_v^0(X_j)^2 \varepsilon_j^2 v(X_i)] + \mathbb{E} [h^{-1} \ell_w^1(X_i, X_i) \ell_w^0(X_i) v(X_i)] \right\}, \\
\mathbb{E} [V_{w,1} E_{v,1}^3 E_{v,2}] &\stackrel{o}{=} 3\mathbb{E} [E_{v,1}^2] \mathbb{E} [V_{w,1} E_{v,1} E_{v,2}], \\
\mathbb{E} [V_{w,2} E_{v,1}^4] &\stackrel{o}{=} 3\mathbb{E} [E_{v,1}^2] \mathbb{E} [V_{w,2} E_{v,1}^2], \\
\mathbb{E} [W_{w,1}^2 E_{v,1}^4] &\stackrel{o}{=} 3\mathbb{E} [E_{v,1}^2] \mathbb{E} [W_{w,1}^2 E_{v,1}^2], \\
\mathbb{E} [V_{w,1}^2 E_{v,1}^4] &\stackrel{o}{=} 3\mathbb{E} [E_{v,1}^2] \mathbb{E} [V_{w,1}^2 E_{v,1}^2].
\end{aligned}$$

The expansion now follows, formally, from the following steps. First, combining the above moments into cumulants. Second, these cumulants may be simplified using that

$$\frac{\sigma_v^2}{\sigma_w^2} = 1 + \mathbb{1}(w \neq v) (\rho^{1+(p+1)} \Omega_{1,\text{bc}} + \rho^{1+2(p+1)} \Omega_{2,\text{bc}})$$

and that in all cases present products such as  $\ell_w^0(X_i)^{k_1} \ell_v^0(X_i)^{k_2}$  and  $\ell_w^1(X_i, X_j)^{k_1} \ell_v^1(X_i, X_j)^{k_2}$  may be replaced with  $\ell_v^0(X_i)^{k_1+k_2}$  and  $\ell_v^1(X_i, X_j)^{k_1+k_2}$ , respectively, provided the arguments match. This is immediate for  $v = w$ , and for  $v \neq w$ , follows because  $\rho \rightarrow 0$  is assumed. This is the analogous step to Eqn. (16) in the density case. For any term of a cumulant with a rate of  $(nh)^{-1}$ ,  $(nh)^{-1/2} \eta_v$ ,  $\eta_v^2$ , or  $\rho^{1+2(p+1)}$  (i.e., the extent of the expansion), these simplifications may be inserted as the remainder will be negligible. Third, with the cumulants in hand, the terms of the expansion are determined as described by e.g., (Hall, 1992a, Chapter 2).

## S.II.7 Complete Simulation Results

In this section we present the results of a simulation study addressing the finite-sample performance of the methods described in the main paper. As with the density estimator, we report empirical coverage probabilities and average interval length of nominal 95% confidence interval for different estimators of a regression functions  $m(x)$  evaluated at values  $x = \{-2/3, -1/3, 0, 1/3, 2/3\}$ . For each replication, the data is generated as i.i.d. draws,  $i = 1, 2, \dots, n$ ,  $n = 500$  as follows:

$$Y = m(x) + \varepsilon, \quad x \sim \mathcal{U}[-1, 1], \quad \varepsilon \sim \mathcal{N}(0, 1)$$

$$\text{Model 1: } m(x) = \sin(4x) + 2 \exp\{-64x^2\}$$

$$\text{Model 2: } m(x) = 2x + 2 \exp\{-64x^2\}$$

$$\text{Model 3: } m(x) = 0.3 \exp\{-4(2x + 1)^2\} + 0.7 \exp\{-16(2x - 1)^2\}$$

$$\text{Model 4: } m(x) = x + 5\phi(10x)$$

$$\text{Model 5: } m(x) = \frac{\sin(3\pi x/2)}{1 + 18x^2[\text{sgn}(x) + 1]}$$

$$\text{Model 6: } m(x) = \frac{\sin(\pi x/2)}{1 + 2x^2[\text{sgn}(x) + 1]}$$

Models 1 to 3 were used by [Fan and Gijbels \(1996\)](#) and [Cattaneo and Farrell \(2013\)](#), while Models 4 to 6 are from [Hall and Horowitz \(2013\)](#), with some originally studied by [Berry et al. \(2002\)](#). The regression functions are plotted in Figure [S.II.1](#) together with the evaluation points used.

We compute confidence intervals for  $m(x)$  using five alternative approaches:

**US:** local-linear estimator using a conventional approach based on undersmoothing ( $I_{\text{us}}$ ).

**Locfit:** local linear estimator computed using default options in the R package `locfit` (see [Loader \(2013\)](#) for implementation details).

**BC:** traditional bias corrected estimator using a local-linear estimator with local-quadratic bias-correction, and  $\rho = 1$  ( $I_{\text{bc}}$ ).

**HH:** local linear estimator using the bootstrapped confidence bands introduced in [Hall and Horowitz \(2013\)](#) (see Remark [10](#) below for additional implementation details).

**RBC:** our proposed local-linear estimator with local-quadratic bias-correction and  $\rho = 1$  using robust standard errors ( $I_{\text{rbc}}$ ).

In all cases the Epanechnikov kernel is used. The bandwidth  $h$  is chosen in three different ways:

- (i) population MSE-optimal choice  $h_{\text{mse}}$ ;
- (ii) estimated ROT optimal coverage error rate  $\hat{h}_{\text{rot}}$ .
- (iii) estimated DPI optimal coverage error rate  $\hat{h}_{\text{dpi}}$ .

For the construction of the variance estimators  $\hat{\sigma}_{\text{us}}^2$  and  $\hat{\sigma}_{\text{rbc}}^2$  we consider HC3 plug-in residuals when forming the  $\Sigma$  matrix. In Table S.II.9 we report empirical coverage and average interval length of RBC 95% Confidence Intervals (only for Model 5) using  $\hat{h}_{\text{mse}}$  for different variance estimators. The results reflect the robustness of the findings to this choice.

The results are presented in detail in the tables and figures below to give a complete picture of the performance of robust bias correction. First, Tables S.II.1-S.II.6 show, for each regression model, respectively, the performance of the five methods above, in terms of empirical coverage and interval length, for all evaluation points and bandwidth choices (recall that  $I_{\text{us}}$  and  $I_{\text{bc}}$  have the same length). Panel A of each shows the coverage and length, while Panel B gives summary statistics for the two fully data-driven bandwidths. Note that in some cases, the population MSE-optimal bandwidth is not defined or is not computable numerically; usually because the bias is too small or other values are too extreme.

The broad conclusion from these tables is that robust bias correction provides excellent coverage and that the data-driven bandwidths perform well and are numerically stable. In almost all cases robust bias correction provides correct coverage, whereas the other methods often, but not always, fail to do so. In cases where there is little to no bias all the methods give good coverage. This can be seen in results for Models 2 and 4, at  $|x| = 2/3$ , far enough away from the “hump” in the center of each, where the true regression function is (nearly) linear. But despite the encouraging results away from the center, only robust bias correction yields good coverage closer to the center ( $|x| = 1/3$ ), when there is more bias. Going further, considering  $x = 0$ , the center of the sharp peak in these models, we see that even robust bias correction fails to provide accurate coverage for  $\hat{h}_{\text{rot}}$ , although  $\hat{h}_{\text{dpi}}$  performs slightly better. At this point, for these models, the bias is too extreme even for robust bias correction to overcome. The results for the other models yield similar lessons.

It is somewhat more difficult to compare interval length using these tables. The comparison is invited for a fixed bandwidth, in which case, by construction, undersmoothing will have a shorter length. However, this ignores the fact that robust bias correction can accommodate

a larger range of bandwidths, and in particular will optimally use a larger bandwidth. For example, robust bias correction has excellent coverage in many cases for  $\hat{h}_{\text{rot}}$ , which is in this case a data-driven MSE-optimal choice (i.e. they coincide). This bandwidth is generally larger than  $\hat{h}_{\text{dpi}}$ , and hence undersmoothing generally covers better with the latter. However, if you compare the length of  $I_{\text{us}}(\hat{h}_{\text{rot}})$  to the length of  $I_{\text{us}}(\hat{h}_{\text{dpi}})$ , we see that robust bias correction compares favorably in terms of length.

Both to better make this point and to illustrate the robustness of  $I_{\text{rbc}}$  to tuning parameter selection, Figures S.II.2–S.II.13 show empirical coverage and length for all six models, and all evaluation points, across a range of bandwidths. The dotted vertical line shows the population MSE-optimal bandwidth (whenever available) for reference. The coverage figures highlight the delicate balance required for undersmoothing to provide correct coverage, and the generally poor performance of traditional bias correction, but show that for a wide range of bandwidths robust bias correction provides correct coverage. Further, interval length is not unduly inflated for bandwidths that provide correct coverage. Again, by construction, undersmoothing will yield shorter intervals for a fixed bandwidth, and this is clear from Figures S.II.8–S.II.13, but it is also clear that robust bias correction can use much larger bandwidths while still maintaining correct coverage.

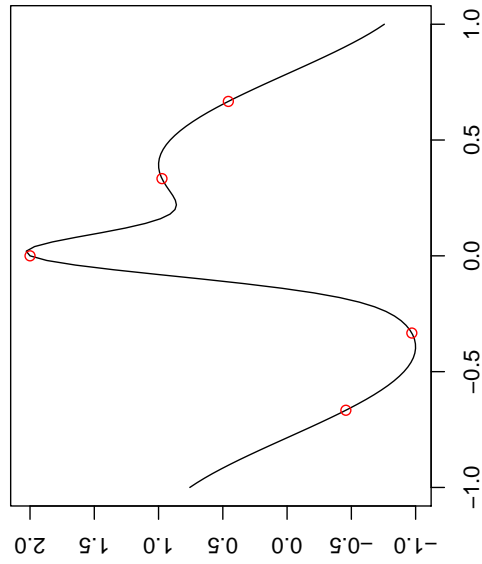
To further illustrate this idea, in Tables S.II.7–S.II.8 we compare average interval length of US and RBC 95% confidence intervals but at different bandwidths. First, in Table S.II.7 we compute average interval length at the largest bandwidth that provides close to correct coverage for each method separately. Note that in all cases these bandwidths are not feasible: these are ex-post findings. Next, in Table S.II.8 we evaluate the performance of US and RBC confidence intervals at certain alternative bandwidths likely to be chosen in practice. First, we evaluate the performance of US confidence intervals at  $h = \lambda \hat{h}_{\text{mse}}$  for  $\lambda = \{0.5; 0.7\}$ . We then compare the performance with RBC confidence intervals computed using the optimal, fully data-driven choices  $\hat{h}_{\text{rot}}$  and  $\hat{h}_{\text{dpi}}$ . Both tables reflect that, once we control for coverage, intervals lengths do not differ systematically between both approaches.

Figures S.II.14–S.II.19 make this same point in a different way. For a range of bandwidths, as in the previous figures, we show the “average position” of  $I_{\text{us}}$  and  $I_{\text{rbc}}$ , where the center of the bar is placed at the average bias and the length of each bar is the average interval length across the simulations. The bars are then color-coded by coverage (green bars having good coverage, fading to red showing undercoverage). These make visually clear that although undersmoothing provides shorter intervals in general, that this comes at the expense of coverage, while robust bias correction provides good coverage for a range of bandwidths, many of which are “large” enough to yield narrow intervals.

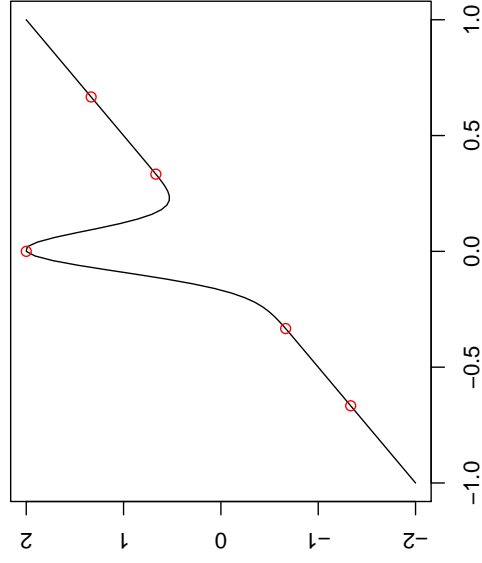
All our methods are implemented in software available from the authors' websites and via the R package `nprobust` available at <https://cran.r-project.org/package=nprobust>.

**Remark 10** (Implementation of Hall and Horowitz (2013)). The column *HH* computes the bootstrapped confidence bands introduced in Hall and Horowitz (2013), following as close as possible their implementation choices. First, we estimate  $m(x)$  using a local linear estimator using the Epanechnikov kernel for our previously discussed bandwidth choices. Standard errors are calculated using their proposed variance estimator  $\hat{\sigma}_{HH}^2 = \kappa \hat{\sigma}^2 / \hat{f}_X(x)$  where  $\kappa = \int K^2$  and  $\hat{f}_X(x)$  is a standard kernel density estimator using a data-driven bandwidth choice  $h_1$ . Then, we use the same estimator for the error variance  $\hat{\sigma}^2 = \sum_{i=1}^n \hat{\varepsilon}_i^2 / n$  and  $\hat{\varepsilon}_i = \tilde{\varepsilon}_i - \bar{\varepsilon}$ ,  $\tilde{\varepsilon}_i = Y_i - \hat{m}(X_i)$ ,  $\bar{\varepsilon} = n^{-1} \sum_{i=1}^n \tilde{\varepsilon}_i$ . Next, we take generate  $B = 500$  bootstrap samples  $\mathcal{Z}^* = \{(X_i, Y_i^*)\}, 1 \leq i \leq n$ , where  $Y_i^* = \hat{m}(X_i) + \varepsilon_i^*$ , with  $\varepsilon_i^*$  obtained by sampling with replacement from the  $\{\hat{\varepsilon}_i\}, 1 \leq i \leq n$ . With these bootstrap samples we can construct the final confidence bands using the adjusted critical values that approximates the estimated coverage error with the selected one. Following their recommendation, the final critical values are taken to be the  $\xi$ -level quantile (for  $\xi = 0.1$ ) obtained by repeating this exercise over a grid of evaluation points, which we choose to be the sequence  $\{x_1, \dots, x_N\} = \{-0.9, -0.8, \dots, 0, \dots, 0.8, 0.9\}$ . ■

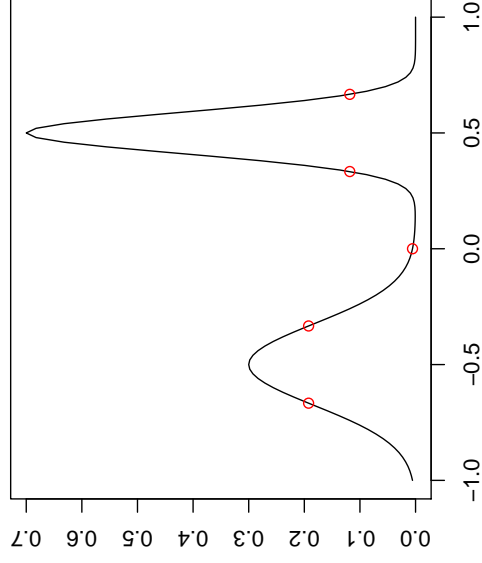
Figure S.II.1: Regression Functions



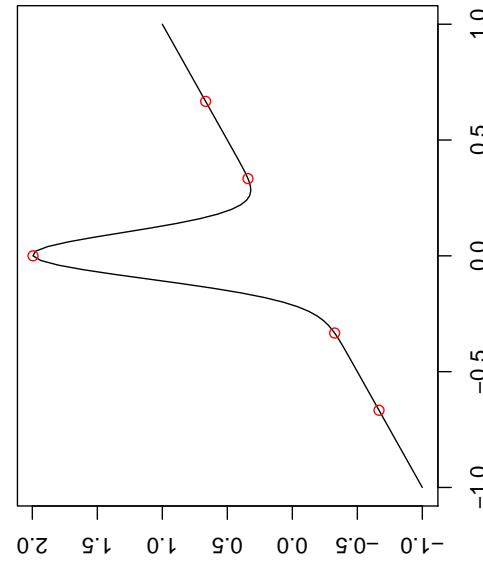
(a) Model 1



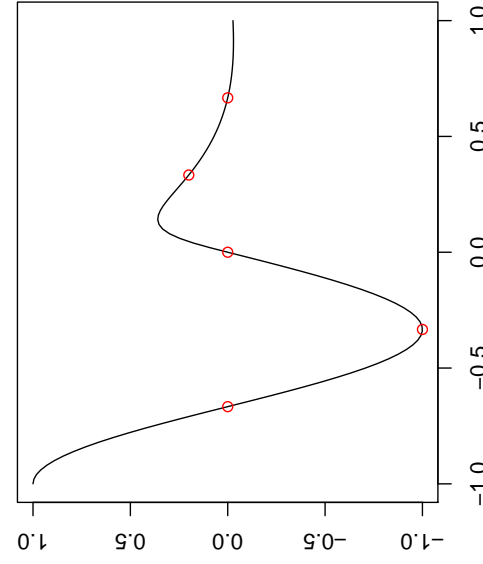
(b) Model 2



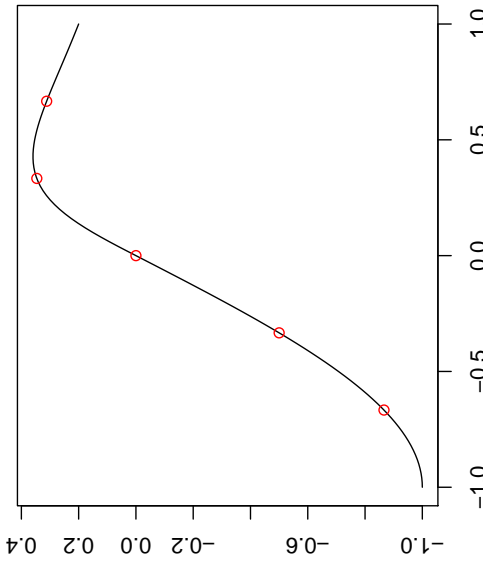
(c) Model 3



(d) Model 4



(e) Model 5



(f) Model 6



Table S.II.1: Simulations Results for Model 1

Panel A: Empirical Coverage and Average Interval Length of 95% Confidence Intervals											
		Bandwidth	Empirical Coverage					Interval Length			
			US	Locfit	BC	HH	RBC	US	Locfit	HH	RBC
$x = -2/3$											
$h_{\text{mse}}$	0.478	57.3	76.8	83.2	31.6	94.3	0.301	0.329	0.197	0.421	
$\hat{h}_{\text{rot}}$	0.202	93.4	94.0	82.9	94.4	94.5	0.439	0.478	0.466	0.628	
$\hat{h}_{\text{dpi}}$	0.132	94.0	94.3	81.7	96.1	93.5	0.575	0.619	0.644	0.818	
$x = -1/3$											
$h_{\text{mse}}$	0.331	4.2	30.7	82.3	1.5	93.1	0.355	0.376	0.276	0.486	
$\hat{h}_{\text{rot}}$	0.486	2.8	9.0	54.0	2.2	63.5	0.326	0.326	0.199	0.417	
$\hat{h}_{\text{dpi}}$	0.284	38.4	61.9	81.6	34.4	92.8	0.385	0.413	0.343	0.538	
$x = 0$											
$h_{\text{mse}}$	0.115	52.4	72.3	83.4	61.4	93.8	0.595	0.623	0.660	0.825	
$\hat{h}_{\text{rot}}$	0.463	0.0	0.0	0.0	0.0	0.0	0.353	0.327	0.199	0.461	
$\hat{h}_{\text{dpi}}$	0.182	5.5	14.8	71.7	7.2	84.6	0.502	0.500	0.504	0.660	
$x = 1/3$											
$h_{\text{mse}}$	0.383	93.2	94.8	77.4	82.3	91.4	0.317	0.353	0.239	0.454	
$\hat{h}_{\text{rot}}$	0.339	94.6	95.4	79.3	87.6	92.6	0.340	0.377	0.281	0.488	
$\hat{h}_{\text{dpi}}$	0.222	93.5	94.0	78.6	93.1	92.2	0.435	0.475	0.450	0.623	
$x = 2/3$											
$h_{\text{mse}}$	0.478	58.8	78.5	83.6	31.4	94.5	0.301	0.330	0.197	0.421	
$\hat{h}_{\text{rot}}$	0.290	88.3	92.4	82.6	82.5	94.4	0.364	0.402	0.323	0.523	
$\hat{h}_{\text{dpi}}$	0.193	91.6	93.0	80.1	91.8	93.2	0.466	0.507	0.502	0.669	

Panel B: Summary Statistics for the Estimated Bandwidths									
	Pop. Par.	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Std. Dev.	
$x = -2/3$									
$\hat{h}_{\text{rot}}$	0.478	0.1573	0.1829	0.1909	0.2018	0.2018	0.6882	0.050	
$\hat{h}_{\text{dpi}}$	-	0.02829	0.1012	0.1232	0.1318	0.1503	0.9401	0.057	
$x = -1/3$									
$\hat{h}_{\text{rot}}$	0.331	0.2201	0.3979	0.4964	0.4855	0.5719	0.7242	0.109	
$\hat{h}_{\text{dpi}}$	-	0.04872	0.2539	0.2862	0.2838	0.3105	1.317	0.079	
$x = 0$									
$\hat{h}_{\text{rot}}$	0.115	0.2886	0.4344	0.4618	0.4635	0.4912	0.6602	0.045	
$\hat{h}_{\text{dpi}}$	-	0.04657	0.1637	0.181	0.1817	0.1994	0.3009	0.028	
$x = 1/3$									
$\hat{h}_{\text{rot}}$	0.383	0.2103	0.2815	0.3353	0.3385	0.3889	0.5925	0.066	
$\hat{h}_{\text{dpi}}$	-	0.02326	0.1666	0.2067	0.2225	0.2626	1.717	0.090	
$x = 2/3$									
$\hat{h}_{\text{rot}}$	0.478	0.212	0.2545	0.281	0.29	0.3189	0.5017	0.044	
$\hat{h}_{\text{dpi}}$	-	0.02727	0.1499	0.1833	0.1934	0.2241	0.9813	0.071	

Notes:

93

(i) US = Undersmoothing, Locfit = R package `locfit` by Loader (2013), BC = Bias Corrected, HH = Hall and Horowitz (2013), RBC = Robust Bias Corrected.

(ii) “Bandwidth” column report the population and average estimated bandwidths choices, as appropriate, for bandwidth  $h$ .

Table S.II.2: Simulations Results for Model 2

Panel A: Empirical Coverage and Average Interval Length of 95% Confidence Intervals										
	Bandwidth	Empirical Coverage					Interval Length			
		US	Locfit	BC	HH	RBC	US	Locfit	HH	RBC
$x = -2/3$										
$h_{\text{mse}}$	-	-	-	-	-	-	-	-	-	-
$\hat{h}_{\text{rot}}$	0.326	95.2	95.3	82.9	86.6	94.8	0.350	0.386	0.280	0.502
$\hat{h}_{\text{dpi}}$	0.219	94.8	95.0	81.9	92.6	94.0	0.443	0.481	0.422	0.632
$x = -1/3$										
$h_{\text{mse}}$	0.706	0.0	0.3	1.0	0.0	3.8	0.253	0.267	0.122	0.355
$\hat{h}_{\text{rot}}$	0.459	0.8	18.7	83.2	0.2	94.2	0.303	0.326	0.189	0.417
$\hat{h}_{\text{dpi}}$	0.311	63.5	77.9	78.6	54.1	91.1	0.362	0.395	0.294	0.514
$x = 0$										
$h_{\text{mse}}$	0.115	52.4	72.4	83.4	49.8	93.8	0.595	0.623	0.573	0.825
$\hat{h}_{\text{rot}}$	0.495	0.0	0.0	0.0	0.0	0.0	0.341	0.315	0.174	0.450
$\hat{h}_{\text{dpi}}$	0.197	1.8	7.6	66.8	1.8	80.3	0.487	0.479	0.432	0.633
$x = 1/3$										
$h_{\text{mse}}$	0.706	0.0	0.3	1.1	0.0	4.8	0.254	0.267	0.122	0.355
$\hat{h}_{\text{rot}}$	0.459	0.7	18.7	84.4	0.0	94.6	0.303	0.326	0.189	0.417
$\hat{h}_{\text{dpi}}$	0.311	63.6	77.3	76.2	54.3	90.2	0.361	0.394	0.294	0.513
$x = 2/3$										
$h_{\text{mse}}$	-	-	-	-	-	-	-	-	-	-
$\hat{h}_{\text{rot}}$	0.323	94.9	95.1	82.4	87.9	94.5	0.351	0.388	0.283	0.505
$\hat{h}_{\text{dpi}}$	0.215	94.4	94.1	80.9	92.7	93.7	0.446	0.487	0.428	0.640

Panel B: Summary Statistics for the Estimated Bandwidths								
	Pop. Par.	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Std. Dev.
$x = -2/3$								
$\hat{h}_{\text{rot}}$	-	0.2046	0.2619	0.2948	0.3262	0.3819	0.5828	NA
$\hat{h}_{\text{dpi}}$	-	0.01928	0.1619	0.2	0.2186	0.2577	1.307	NA
$x = -1/3$								
$\hat{h}_{\text{rot}}$	0.706	0.3176	0.4301	0.4558	0.4589	0.4843	0.6388	0.042
$\hat{h}_{\text{dpi}}$	-	0.1275	0.2603	0.2927	0.3106	0.3373	1.718	0.090
$x = 0$								
$\hat{h}_{\text{rot}}$	0.115	0.3844	0.4694	0.4921	0.4949	0.5171	0.6434	0.036
$\hat{h}_{\text{dpi}}$	-	0.0537	0.18	0.1957	0.1968	0.2125	0.2953	0.025
$x = 1/3$								
$\hat{h}_{\text{rot}}$	0.706	0.2997	0.4299	0.4557	0.459	0.4846	0.6554	0.042
$\hat{h}_{\text{dpi}}$	-	0.1289	0.2602	0.2922	0.3105	0.3392	1.931	0.093
$x = 2/3$								
$\hat{h}_{\text{rot}}$	-	0.2065	0.2596	0.29	0.3226	0.3762	0.5824	0.082
$\hat{h}_{\text{dpi}}$	-	0.02684	0.1578	0.1957	0.2146	0.2558	1.081	0.089

Notes:

94

(i) US = Undersmoothing, Locfit = R package `locfit` by Loader (2013), BC = Bias Corrected, HH = Hall and Horowitz (2013), RBC = Robust Bias Corrected.

(ii) “Bandwidth” column report the population and average estimated bandwidths choices, as appropriate, for bandwidth  $h$ .

Table S.II.3: Simulations Results for Model 3

Panel A: Empirical Coverage and Average Interval Length of 95% Confidence Intervals										
	Bandwidth	Empirical Coverage					Interval Length			
		US	Locfit	BC	HH	RBC	US	Locfit	HH	RBC
$x = -2/3$										
$h_{\text{mse}}$	-	-	-	-	-	-	-	-	-	-
$\hat{h}_{\text{rot}}$	0.532	91.2	92.1	85.7	64.3	95.2	0.299	0.313	0.166	0.405
$\hat{h}_{\text{dpi}}$	0.346	93.3	93.5	82.5	81.1	94.1	0.346	0.374	0.257	0.495
$x = -1/3$										
$h_{\text{mse}}$	-	-	-	-	-	-	-	-	-	-
$\hat{h}_{\text{rot}}$	0.696	80.1	87.0	82.1	43.2	94.8	0.232	0.253	0.116	0.335
$\hat{h}_{\text{dpi}}$	0.491	90.7	92.6	81.0	68.2	94.0	0.281	0.307	0.173	0.405
$x = 0$										
$h_{\text{mse}}$	0.976	13.5	19.7	39.8	1.3	61.9	0.198	0.214	0.082	0.283
$\hat{h}_{\text{rot}}$	0.696	34.3	63.2	84.9	7.7	95.7	0.234	0.253	0.116	0.333
$\hat{h}_{\text{dpi}}$	0.491	79.9	87.7	79.0	56.1	92.3	0.282	0.308	0.174	0.406
$x = 1/3$										
$h_{\text{mse}}$	0.246	77.8	86.1	79.6	67.3	92.7	0.393	0.423	0.326	0.563
$\hat{h}_{\text{rot}}$	0.695	86.6	83.2	49.7	52.7	71.8	0.237	0.253	0.116	0.343
$\hat{h}_{\text{dpi}}$	0.494	76.5	71.5	52.5	47.2	73.3	0.285	0.307	0.172	0.410
$x = 2/3$										
$h_{\text{mse}}$	0.246	79.0	85.7	79.7	68.2	92.8	0.393	0.424	0.327	0.564
$\hat{h}_{\text{rot}}$	0.505	78.4	75.7	46.3	47.5	69.1	0.307	0.320	0.176	0.422
$\hat{h}_{\text{dpi}}$	0.325	78.3	82.7	70.8	60.2	88.0	0.360	0.387	0.274	0.516

Panel B: Summary Statistics for the Estimated Bandwidths								
	Pop. Par.	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Std. Dev.
$x = -2/3$								
$\hat{h}_{\text{rot}}$	-	0.2424	0.4328	0.5297	0.5317	0.6234	0.855	0.122
$\hat{h}_{\text{dpi}}$	-	0.07286	0.269	0.3288	0.3458	0.3989	1.65	0.119
$x = -1/3$								
$\hat{h}_{\text{rot}}$	-	0.4494	0.6663	0.7015	0.6957	0.7301	0.8429	0.049
$\hat{h}_{\text{dpi}}$	-	0.2405	0.406	0.4573	0.4913	0.5358	2.668	0.140
$x = 0$								
$\hat{h}_{\text{rot}}$	0.976	0.4981	0.6655	0.7024	0.696	0.7324	0.8301	0.051
$\hat{h}_{\text{dpi}}$	-	0.2371	0.4036	0.4592	0.4913	0.5368	2.671	0.147
$x = 1/3$								
$\hat{h}_{\text{rot}}$	0.246	0.4885	0.6638	0.702	0.6953	0.7321	0.8256	0.052
$\hat{h}_{\text{dpi}}$	-	0.2426	0.4062	0.4596	0.4942	0.5372	2.992	0.153
$x = 2/3$								
$\hat{h}_{\text{rot}}$	0.246	0.204	0.3964	0.4995	0.5049	0.6076	0.8263	0.131
$\hat{h}_{\text{dpi}}$	-	0.06684	0.24	0.3105	0.3252	0.381	1.611	0.118

Notes:

95

(i) US = Undersmoothing, Locfit = R package `locfit` by Loader (2013), BC = Bias Corrected, HH = Hall and Horowitz (2013), RBC = Robust Bias Corrected.

(ii) “Bandwidth” column report the population and average estimated bandwidths choices, as appropriate, for bandwidth  $h$ .

Table S.II.4: Simulations Results for Model 4

Panel A: Empirical Coverage and Average Interval Length of 95% Confidence Intervals										
	Bandwidth	Empirical Coverage					Interval Length			
		US	Locfit	BC	HH	RBC	US	Locfit	HH	RBC
$x = -2/3$										
$h_{\text{mse}}$	-	-	-	-	-	-	-	-	-	-
$\hat{h}_{\text{rot}}$	0.310	95.1	95.3	83.2	88.3	95.0	0.357	0.392	0.294	0.512
$\hat{h}_{\text{dpi}}$	0.208	94.5	95.0	81.9	93.3	93.9	0.452	0.491	0.437	0.646
$x = -1/3$										
$h_{\text{mse}}$	0.466	0.3	8.1	76.8	0.0	90.3	0.300	0.322	0.184	0.412
$\hat{h}_{\text{rot}}$	0.439	0.6	14.6	82.5	0.1	94.1	0.308	0.332	0.198	0.425
$\hat{h}_{\text{dpi}}$	0.304	56.5	73.8	79.5	47.9	91.2	0.366	0.398	0.301	0.519
$x = 0$										
$h_{\text{mse}}$	0.127	51.8	71.9	83.4	50.8	93.9	0.564	0.592	0.556	0.784
$\hat{h}_{\text{rot}}$	0.472	0.0	0.0	0.0	0.0	0.0	0.348	0.320	0.183	0.446
$\hat{h}_{\text{dpi}}$	0.188	6.6	19.5	76.1	6.6	88.8	0.483	0.489	0.449	0.645
$x = 1/3$										
$h_{\text{mse}}$	0.466	0.2	8.9	75.5	0.0	89.7	0.300	0.321	0.184	0.411
$\hat{h}_{\text{rot}}$	0.439	0.4	14.5	82.9	0.0	94.1	0.308	0.331	0.197	0.425
$\hat{h}_{\text{dpi}}$	0.304	57.3	73.3	77.0	48.8	90.5	0.366	0.397	0.301	0.519
$x = 2/3$										
$h_{\text{mse}}$	-	-	-	-	-	-	-	-	-	-
$\hat{h}_{\text{rot}}$	0.307	94.8	94.9	82.2	88.7	94.5	0.358	0.395	0.296	0.516
$\hat{h}_{\text{dpi}}$	0.204	94.0	94.2	81.3	92.6	93.9	0.458	0.498	0.444	0.657

Panel B: Summary Statistics for the Estimated Bandwidths								
	Pop. Par.	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Std. Dev.
$x = -2/3$								
$\hat{h}_{\text{rot}}$	-	0.2014	0.2547	0.2807	0.3101	0.3445	0.5721	0.078
$\hat{h}_{\text{dpi}}$	-	0.02144	0.1547	0.1907	0.2076	0.2423	1.878	0.090
$x = -1/3$								
$\hat{h}_{\text{rot}}$	0.466	0.3059	0.4129	0.4367	0.439	0.4624	0.5956	0.038
$\hat{h}_{\text{dpi}}$	-	0.1261	0.253	0.2857	0.3038	0.3325	1.192	0.087
$x = 0$								
$\hat{h}_{\text{rot}}$	0.127	0.3694	0.4497	0.4703	0.4722	0.4923	0.6034	0.032
$\hat{h}_{\text{dpi}}$	-	0.07306	0.1725	0.1869	0.188	0.2022	0.317	0.023
$x = 1/3$								
$\hat{h}_{\text{rot}}$	0.466	0.2918	0.413	0.4367	0.439	0.4626	0.6154	0.038
$\hat{h}_{\text{dpi}}$	-	0.1294	0.2527	0.2853	0.3043	0.3309	1.856	0.094
$x = 2/3$								
$\hat{h}_{\text{rot}}$	-	0.2032	0.2527	0.2776	0.3069	0.34	0.5755	0.076
$\hat{h}_{\text{dpi}}$	-	0.02984	0.1513	0.1869	0.2041	0.2418	1.03	0.084

Notes:

96

(i) US = Undersmoothing, Locfit = R package `locfit` by Loader (2013), BC = Bias Corrected, HH = Hall and Horowitz (2013), RBC = Robust Bias Corrected.

(ii) “Bandwidth” column report the population and average estimated bandwidths choices, as appropriate, for bandwidth  $h$ .

Table S.II.5: Simulations Results for Model 5

Panel A: Empirical Coverage and Average Interval Length of 95% Confidence Intervals										
	Bandwidth	Empirical Coverage					Interval Length			
		US	Locfit	BC	HH	RBC	US	Locfit	HH	RBC
$x = -2/3$										
$h_{\text{mse}}$	-	-	-	-	-	-	-	-	-	-
$\hat{h}_{\text{rot}}$	0.251	95.1	94.7	82.9	90.8	94.8	0.390	0.423	0.338	0.560
$\hat{h}_{\text{dpi}}$	0.166	94.8	94.4	81.8	93.5	93.7	0.505	0.544	0.479	0.722
$x = -1/3$										
$h_{\text{mse}}$	0.307	43.0	69.0	84.0	25.3	94.5	0.354	0.379	0.271	0.502
$\hat{h}_{\text{rot}}$	0.405	9.7	27.7	81.9	4.9	93.4	0.316	0.333	0.209	0.439
$\hat{h}_{\text{dpi}}$	0.283	56.5	70.7	80.6	48.2	92.8	0.380	0.409	0.316	0.540
$x = 0$										
$h_{\text{mse}}$	-	-	-	-	-	-	-	-	-	-
$\hat{h}_{\text{rot}}$	0.475	24.8	50.4	79.4	5.5	92.8	0.286	0.308	0.176	0.409
$\hat{h}_{\text{dpi}}$	0.318	74.4	83.7	80.3	61.1	92.6	0.354	0.383	0.279	0.507
$x = 1/3$										
$h_{\text{mse}}$	0.821	3.3	37.3	81.2	0.1	93.5	0.226	0.240	0.102	0.318
$\hat{h}_{\text{rot}}$	0.536	72.1	88.1	77.3	43.9	92.1	0.268	0.292	0.158	0.384
$\hat{h}_{\text{dpi}}$	0.370	89.9	92.1	78.5	78.4	92.9	0.327	0.356	0.241	0.470
$x = 2/3$										
$h_{\text{mse}}$	0.886	91.0	94.2	74.8	46.0	79.9	0.288	0.312	0.107	0.315
$\hat{h}_{\text{rot}}$	0.400	93.5	93.9	82.7	79.5	94.4	0.318	0.341	0.218	0.453
$\hat{h}_{\text{dpi}}$	0.265	93.9	93.9	81.3	88.4	93.6	0.391	0.425	0.339	0.562

Panel B: Summary Statistics for the Estimated Bandwidths								
	Pop. Par.	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Std. Dev.
$x = -2/3$								
$\hat{h}_{\text{rot}}$	-	0.1888	0.2253	0.2396	0.2513	0.262	0.5376	0.044
$\hat{h}_{\text{dpi}}$	-	0.02597	0.1302	0.1563	0.1665	0.1929	0.8877	0.065
$x = -1/3$								
$\hat{h}_{\text{rot}}$	0.307	0.245	0.3677	0.4093	0.4045	0.4421	0.5661	0.053
$\hat{h}_{\text{dpi}}$	-	0.04247	0.2272	0.2666	0.2825	0.317	1.897	0.093
$x = 0$								
$\hat{h}_{\text{rot}}$	-	0.362	0.4439	0.4707	0.4747	0.5004	0.6879	0.044
$\hat{h}_{\text{dpi}}$	-	0.157	0.259	0.2983	0.318	0.3538	1.57	0.096
$x = 1/3$								
$\hat{h}_{\text{rot}}$	0.821	0.2897	0.479	0.5284	0.5364	0.5959	0.7665	0.079
$\hat{h}_{\text{dpi}}$	-	0.1258	0.2967	0.3465	0.3699	0.4137	1.508	0.115
$x = 2/3$								
$\hat{h}_{\text{rot}}$	0.886	0.2545	0.3442	0.376	0.3998	0.4205	0.7831	0.086
$\hat{h}_{\text{dpi}}$	-	0.06169	0.2045	0.2436	0.2651	0.3033	0.8067	0.090

Notes:

97

(i) US = Undersmoothing, Locfit = R package `locfit` by Loader (2013), BC = Bias Corrected, HH = Hall and Horowitz (2013), RBC = Robust Bias Corrected.

(ii) “Bandwidth” column report the population and average estimated bandwidths choices, as appropriate, for bandwidth  $h$ .

Table S.II.6: Simulations Results for Model 6

Panel A: Empirical Coverage and Average Interval Length of 95% Confidence Intervals											
		Bandwidth	Empirical Coverage					Interval Length			
			US	Locfit	BC	HH	RBC	US	Locfit	HH	RBC
$x = -2/3$											
$h_{\text{mse}}$	0.782	88.9	89.4	90.6	45.7	94.6	0.288	0.298	0.113	0.332	
$\hat{h}_{\text{rot}}$	0.565	90.5	91.4	85.7	60.8	94.7	0.294	0.302	0.152	0.391	
$\hat{h}_{\text{dpi}}$	0.371	93.1	93.2	81.7	77.4	93.9	0.333	0.358	0.233	0.475	
$x = -1/3$											
$h_{\text{mse}}$	0.975	80.1	83.2	77.2	34.3	91.1	0.210	0.217	0.084	0.295	
$\hat{h}_{\text{rot}}$	0.578	91.5	93.4	83.6	64.3	95.2	0.254	0.276	0.139	0.366	
$\hat{h}_{\text{dpi}}$	0.411	93.8	94.0	82.4	78.2	94.0	0.309	0.336	0.207	0.445	
$x = 0$											
$h_{\text{mse}}$	-	-	-	-	-	-	-	-	-	-	
$\hat{h}_{\text{rot}}$	0.562	87.0	91.1	81.6	60.4	94.6	0.258	0.280	0.142	0.372	
$\hat{h}_{\text{dpi}}$	0.401	90.9	92.2	80.5	74.6	93.3	0.312	0.340	0.212	0.450	
$x = 1/3$											
$h_{\text{mse}}$	0.616	51.9	73.4	81.8	18.5	94.7	0.246	0.266	0.129	0.353	
$\hat{h}_{\text{rot}}$	0.546	66.6	78.9	80.9	36.5	93.7	0.262	0.284	0.147	0.376	
$\hat{h}_{\text{dpi}}$	0.389	83.3	87.3	79.6	66.7	93.3	0.318	0.345	0.219	0.458	
$x = 2/3$											
$h_{\text{mse}}$	-	-	-	-	-	-	-	-	-	-	
$\hat{h}_{\text{rot}}$	0.462	94.9	94.5	83.3	75.0	94.4	0.303	0.317	0.180	0.427	
$\hat{h}_{\text{dpi}}$	0.307	94.4	94.2	81.3	84.8	93.7	0.362	0.392	0.279	0.520	

Panel B: Summary Statistics for the Estimated Bandwidths								
	Pop. Par.	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Std. Dev.
$x = -2/3$								
$\hat{h}_{\text{rot}}$	0.782	0.2668	0.5164	0.5764	0.5647	0.6229	0.7842	0.084
$\hat{h}_{\text{dpi}}$	-	0.1125	0.2998	0.3481	0.3707	0.4131	2.057	0.121
$x = -1/3$								
$\hat{h}_{\text{rot}}$	0.975	0.4066	0.5321	0.5717	0.5778	0.6184	0.7935	0.063
$\hat{h}_{\text{dpi}}$	-	0.1991	0.3296	0.382	0.4106	0.4524	2.029	0.129
$x = 0$								
$\hat{h}_{\text{rot}}$	-	0.4028	0.5188	0.5565	0.5622	0.5992	0.7875	0.059
$\hat{h}_{\text{dpi}}$	-	0.1903	0.3237	0.3729	0.4009	0.445	2.431	0.125
$x = 1/3$								
$\hat{h}_{\text{rot}}$	0.616	0.3523	0.5043	0.5402	0.546	0.582	0.7959	0.058
$\hat{h}_{\text{dpi}}$	-	0.166	0.3113	0.3609	0.3892	0.4289	2.203	0.131
$x = 2/3$								
$\hat{h}_{\text{rot}}$	-	0.262	0.4121	0.4535	0.4618	0.5018	0.8093	0.069
$\hat{h}_{\text{dpi}}$	-	0.1082	0.2444	0.2868	0.3071	0.3461	1.515	0.100

**Notes:**

(i) US = Undersmoothing, Locfit = R package `locfit` by Loader (2013), BC = Bias Corrected, HH = Hall and Horowitz (2013), RBC = Robust Bias Corrected.

(ii) “Bandwidth” column report the population and average estimated bandwidths choices, as appropriate, for bandwidth  $h_n$ .

Table S.II.7: Empirical Coverage and Average Interval Length of 95% Confidence Intervals

	US			RBC		
	$h$	EC	IL	$h$	EC	IL
<b>Model 1</b>						
$x = -2/3$	0.140	94.8	0.523	0.420	94.8	0.442
$x = -1/3$	0.100	94.7	0.625	0.420	94.8	0.434
$x = 0$	0.100	71.3	0.640	0.100	93.7	0.893
$x = 1/3$	0.300	94.6	0.355	0.440	94.3	0.425
$x = 2/3$	0.100	95.0	0.624	0.260	94.9	0.546
<b>Model 2</b>						
$x = -2/3$	0.180	94.9	0.459	0.540	94.9	0.399
$x = -1/3$	0.140	94.8	0.524	0.440	94.9	0.424
$x = 0$	0.100	71.3	0.640	0.100	93.7	0.893
$x = 1/3$	0.140	94.5	0.522	0.440	94.2	0.424
$x = 2/3$	0.260	94.9	0.380	0.280	94.9	0.525
<b>Model 3</b>						
$x = -2/3$	0.140	94.9	0.523	0.420	94.9	0.442
$x = -1/3$	0.200	94.9	0.435	0.400	94.9	0.440
$x = 0$	0.100	94.7	0.628	0.680	94.7	0.337
$x = 1/3$	0.100	93.9	0.623	0.100	94.0	0.887
$x = 2/3$	0.100	94.6	0.624	0.180	94.9	0.658
<b>Model 4</b>						
$x = -2/3$	0.180	94.9	0.459	0.520	94.8	0.406
$x = -1/3$	0.100	94.8	0.625	0.400	94.8	0.444
$x = 0$	0.100	79.3	0.636	0.100	93.9	0.893
$x = 1/3$	0.100	94.4	0.623	0.400	94.2	0.443
$x = 2/3$	0.320	94.9	0.342	0.280	94.9	0.525
<b>Model 5</b>						
$x = -2/3$	0.180	94.9	0.459	0.200	94.8	0.624
$x = -1/3$	0.100	94.7	0.625	0.180	94.6	0.658
$x = 0$	0.100	94.6	0.628	0.240	94.4	0.572
$x = 1/3$	0.140	94.6	0.522	0.260	94.3	0.545
$x = 2/3$	0.200	94.8	0.434	0.280	94.9	0.525
<b>Model 6</b>						
$x = -2/3$	0.140	94.9	0.523	0.600	94.9	0.379
$x = -1/3$	0.140	94.8	0.524	0.420	94.9	0.429
$x = 0$	0.100	94.8	0.628	0.600	94.9	0.359
$x = 1/3$	0.140	94.5	0.522	0.480	94.4	0.401
$x = 2/3$	0.260	94.8	0.380	0.420	94.9	0.442

**Notes:** Bandwidths are selected ex post as the largest bandwidths yielding good coverage, and as can not be made feasible

Table S.II.8: Empirical Coverage and Average Interval Length of 95% Confidence Intervals

	US ( $\lambda = 0.5$ )		US ( $\lambda = 0.7$ )		RBC ( $\hat{h}_{\text{rot}}$ )		RBC ( $\hat{h}_{\text{dpi}}$ )	
	EC	IL	EC	IL	EC	IL	EC	IL
<b>Model 1</b>								
$x = -2/3$	94.4	0.630	94.7	0.528	94.3	0.630	93.5	0.818
$x = -1/3$	56.5	0.410	21.1	0.362	63.3	0.417	92.8	0.538
$x = 0$	0.0	0.466	0.0	0.414	0.0	0.463	84.6	0.660
$x = 1/3$	93.5	0.479	94.1	0.404	92.4	0.486	92.2	0.623
$x = 2/3$	95.0	0.519	93.3	0.436	94.9	0.522	93.2	0.669
<b>Model 2</b>								
$x = -2/3$	94.9	0.495	95.2	0.416	95.1	0.503	94.0	0.632
$x = -1/3$	92.7	0.408	57.9	0.350	94.4	0.417	91.1	0.514
$x = 0$	0.0	0.455	0.0	0.403	0.0	0.451	80.3	0.633
$x = 1/3$	92.4	0.407	58.0	0.350	93.9	0.417	90.2	0.513
$x = 2/3$	95.3	0.496	95.0	0.417	94.9	0.503	93.7	0.640
<b>Model 3</b>								
$x = -2/3$	94.4	0.384	93.9	0.329	94.9	0.405	94.1	0.495
$x = -1/3$	93.9	0.328	91.4	0.277	94.1	0.336	94.0	0.405
$x = 0$	94.5	0.329	87.5	0.277	95.8	0.334	92.3	0.406
$x = 1/3$	71.2	0.331	77.5	0.281	73.0	0.343	73.3	0.410
$x = 2/3$	81.4	0.399	74.7	0.343	68.9	0.423	88.0	0.516
<b>Model 4</b>								
$x = -2/3$	94.9	0.507	95.1	0.426	95.0	0.513	93.9	0.646
$x = -1/3$	90.2	0.418	51.8	0.358	93.9	0.425	91.2	0.519
$x = 0$	0.0	0.451	0.0	0.403	0.0	0.448	88.8	0.645
$x = 1/3$	90.3	0.417	52.3	0.357	93.5	0.424	90.5	0.519
$x = 2/3$	95.4	0.508	95.0	0.427	94.9	0.514	93.9	0.657
<b>Model 5</b>								
$x = -2/3$	94.6	0.560	95.0	0.470	94.4	0.562	93.7	0.722
$x = -1/3$	85.1	0.437	55.0	0.370	93.1	0.440	92.8	0.540
$x = 0$	90.8	0.402	73.5	0.340	92.0	0.410	92.6	0.507
$x = 1/3$	94.4	0.378	94.1	0.319	92.2	0.385	92.9	0.470
$x = 2/3$	95.2	0.442	94.7	0.373	95.0	0.454	93.6	0.562
<b>Model 6</b>								
$x = -2/3$	94.3	0.368	93.2	0.317	94.9	0.392	93.9	0.475
$x = -1/3$	94.9	0.362	94.4	0.305	94.5	0.366	94.0	0.445
$x = 0$	94.1	0.367	93.0	0.309	94.9	0.372	93.3	0.450
$x = 1/3$	92.6	0.372	86.8	0.313	93.6	0.377	93.3	0.458
$x = 2/3$	94.8	0.407	94.5	0.344	94.7	0.427	93.7	0.520

**Notes:** Undersmoothing is implemented using bandwidths  $h = \lambda \hat{h}_{\text{mse}}$  for  $\lambda = \{0.5; 0.7\}$ , in the columns labeled as such.



Table S.II.9: Empirical Coverage and Average Interval Length of RBC 95% Confidence Intervals for Model 5, for Different Variance Estimators

	$h$	EC	IL
$x = -2/3$			
$HC_0$	0.248	94.2	0.561
$HC_1$	0.248	94.5	0.556
$HC_2$	0.249	94.5	0.562
$HC_3$	0.249	94.5	0.559
$NN$	0.251	94.8	0.560
$x = -1/3$			
$HC_0$	0.400	92.3	0.440
$HC_1$	0.403	91.9	0.436
$HC_2$	0.404	92.0	0.439
$HC_3$	0.404	91.9	0.437
$NN$	0.405	93.4	0.439
$x = 0$			
$HC_0$	0.472	92.2	0.403
$HC_1$	0.471	92.2	0.407
$HC_2$	0.472	92.6	0.409
$HC_3$	0.472	92.4	0.408
$NN$	0.475	92.8	0.409
$x = 1/3$			
$HC_0$	0.543	90.6	0.378
$HC_1$	0.535	91.0	0.382
$HC_2$	0.536	91.0	0.384
$HC_3$	0.536	91.0	0.383
$NN$	0.536	92.1	0.384
$x = 2/3$			
$HC_0$	0.403	93.6	0.451
$HC_1$	0.399	93.8	0.449
$HC_2$	0.400	94.1	0.452
$HC_3$	0.400	94.0	0.451
$NN$	0.400	94.4	0.453

**Notes:**

- (i) The  $h$  column reports the average estimated bandwidths  $\hat{h}_{\text{rot}}$ .

Figure S.II.2: Empirical Coverage of 95% Confidence Intervals - Model 1

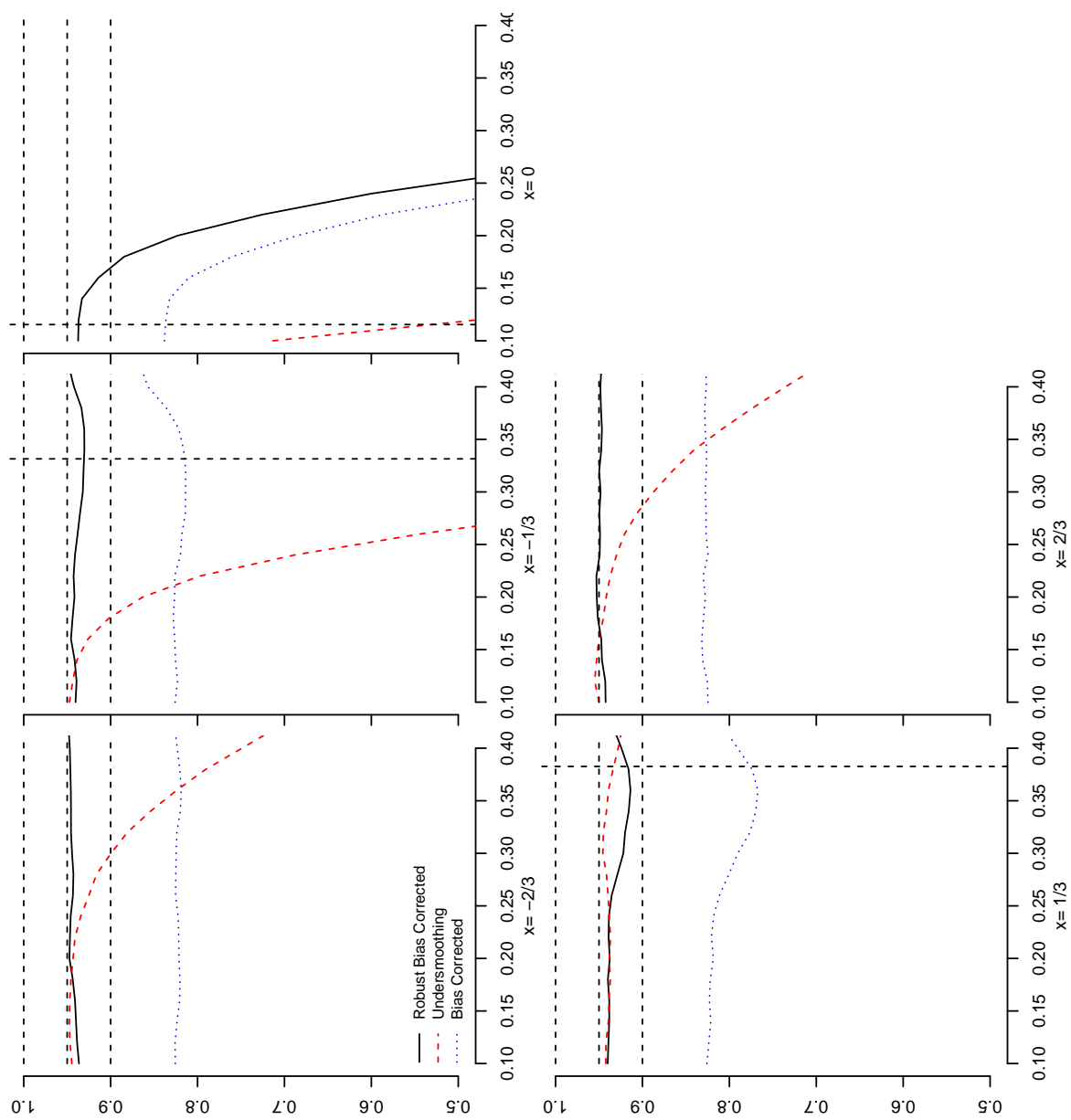


Figure S.II.3: Empirical Coverage of 95% Confidence Intervals - Model 2

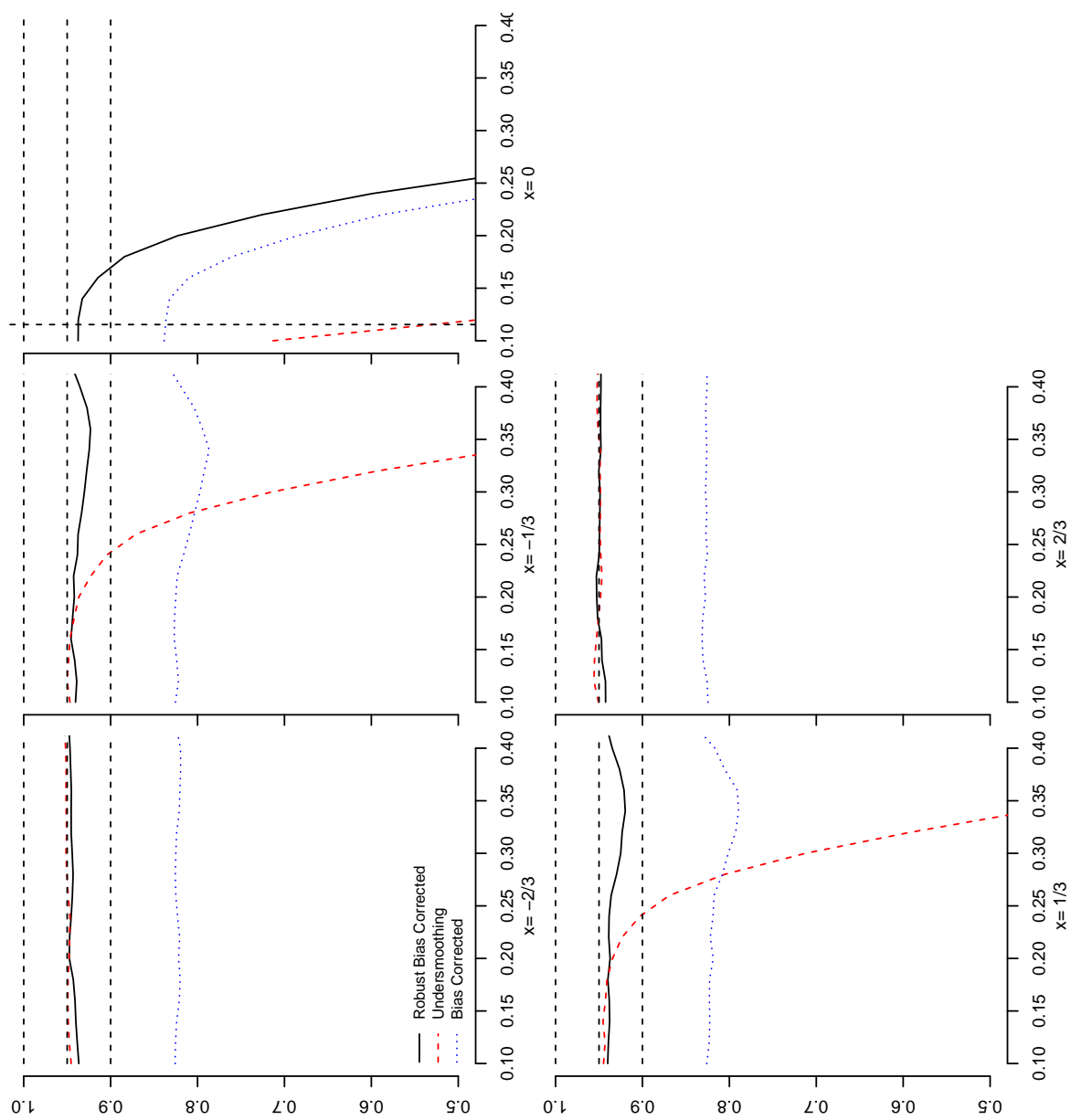


Figure S.II.4: Empirical Coverage of 95% Confidence Intervals - Model 3

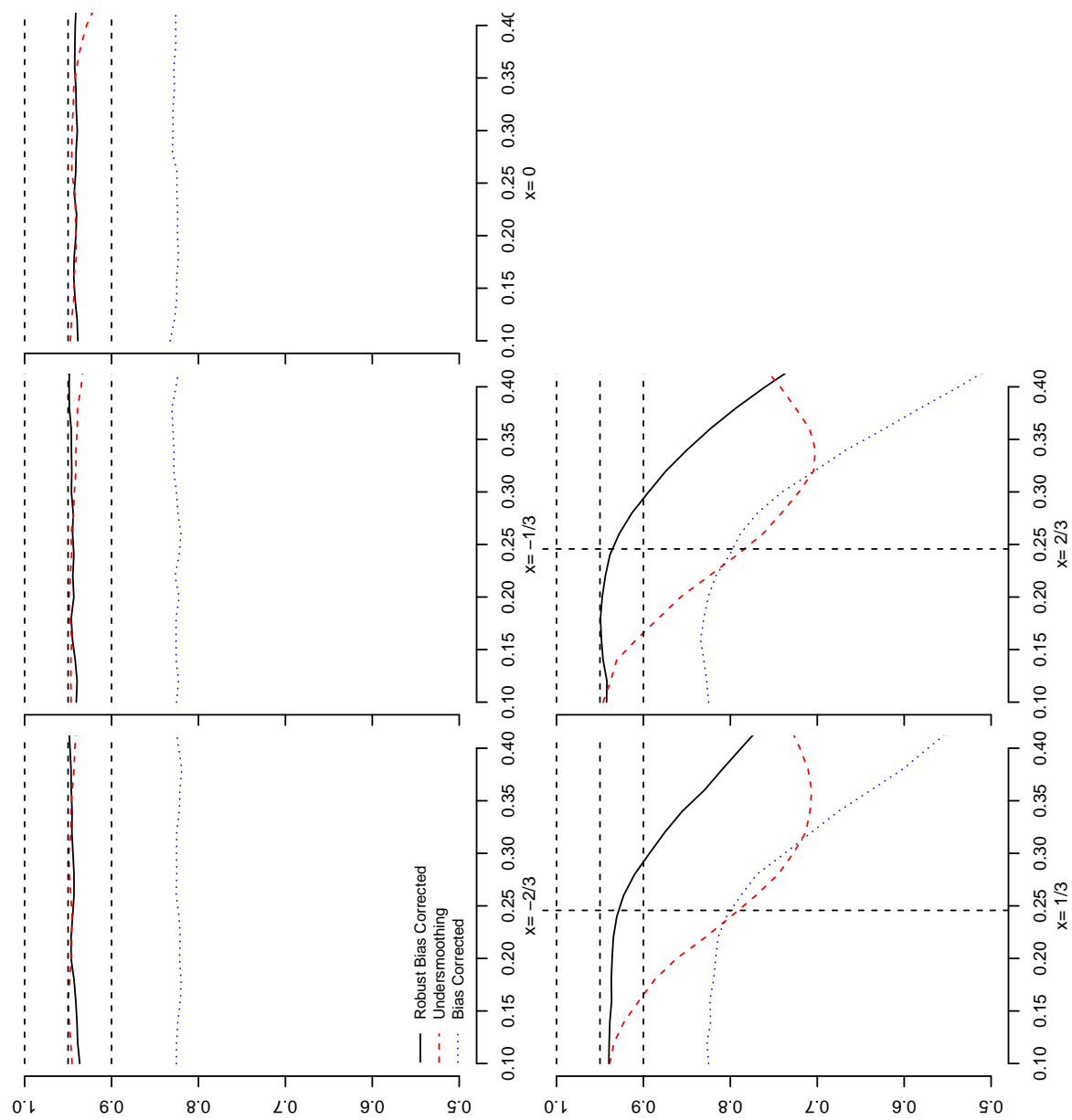


Figure S.II.5: Empirical Coverage of 95% Confidence Intervals - Model 4

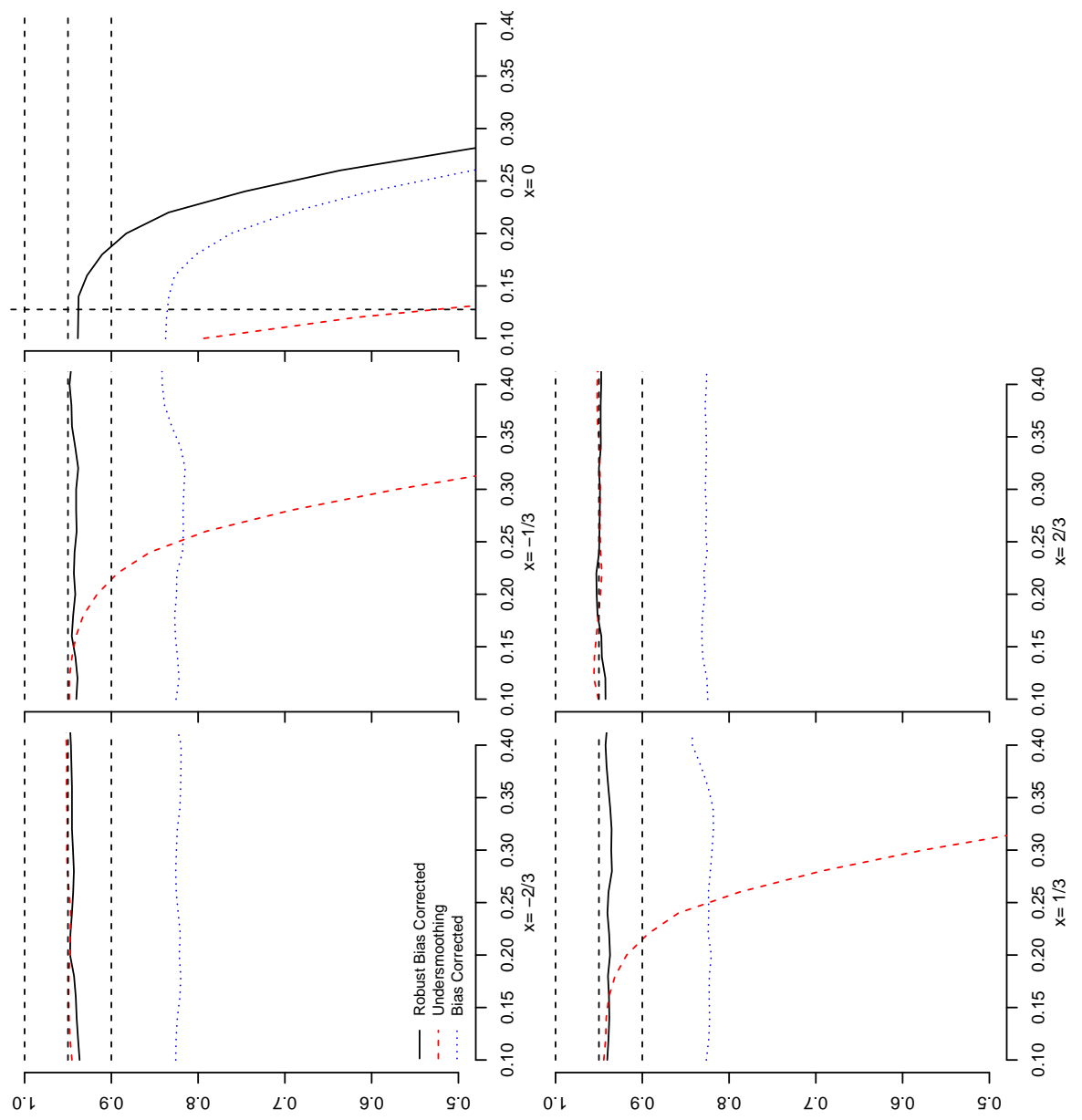


Figure S.II.6: Empirical Coverage of 95% Confidence Intervals - Model 5

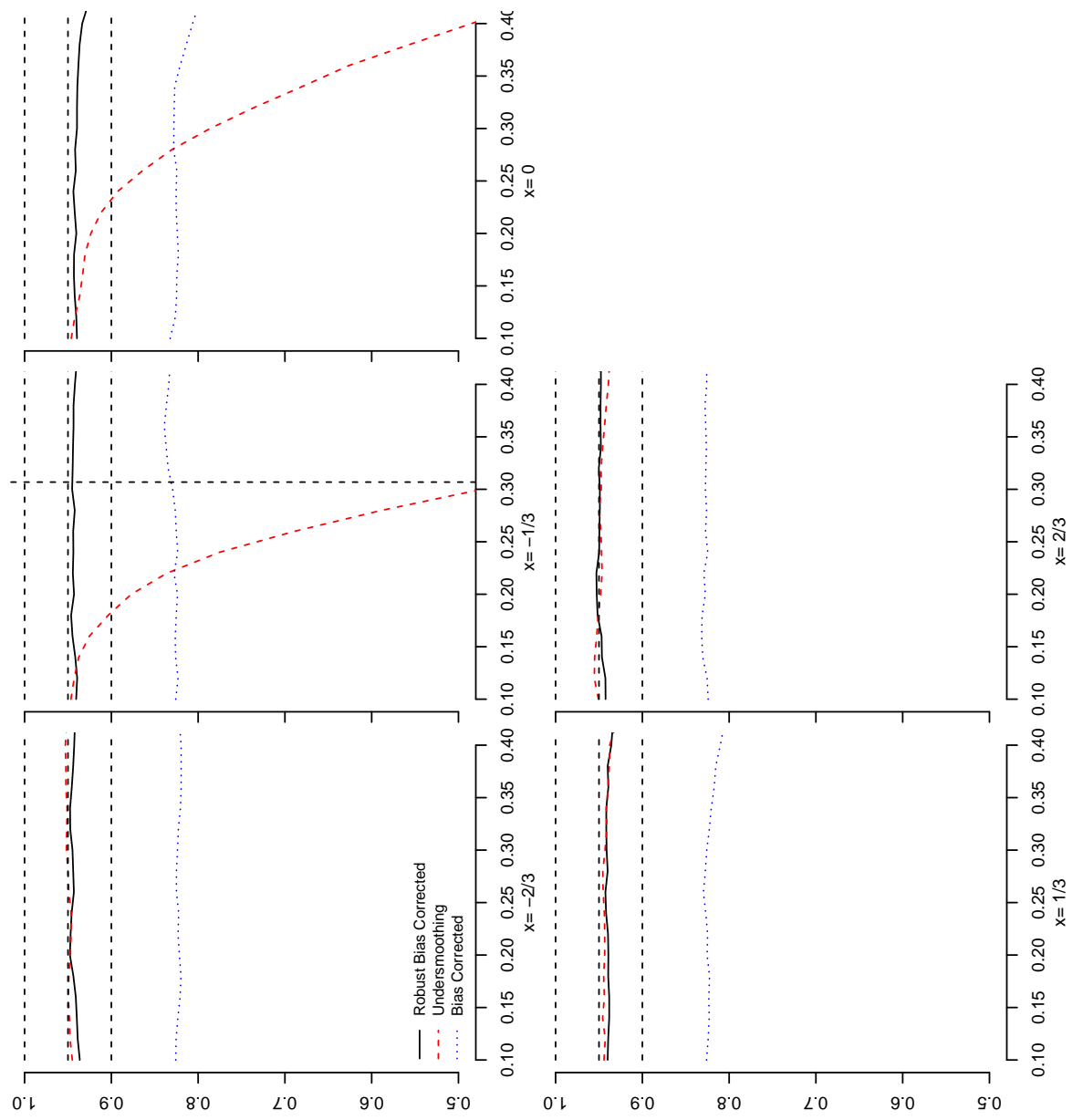


Figure S.II.7: Empirical Coverage of 95% Confidence Intervals - Model 6

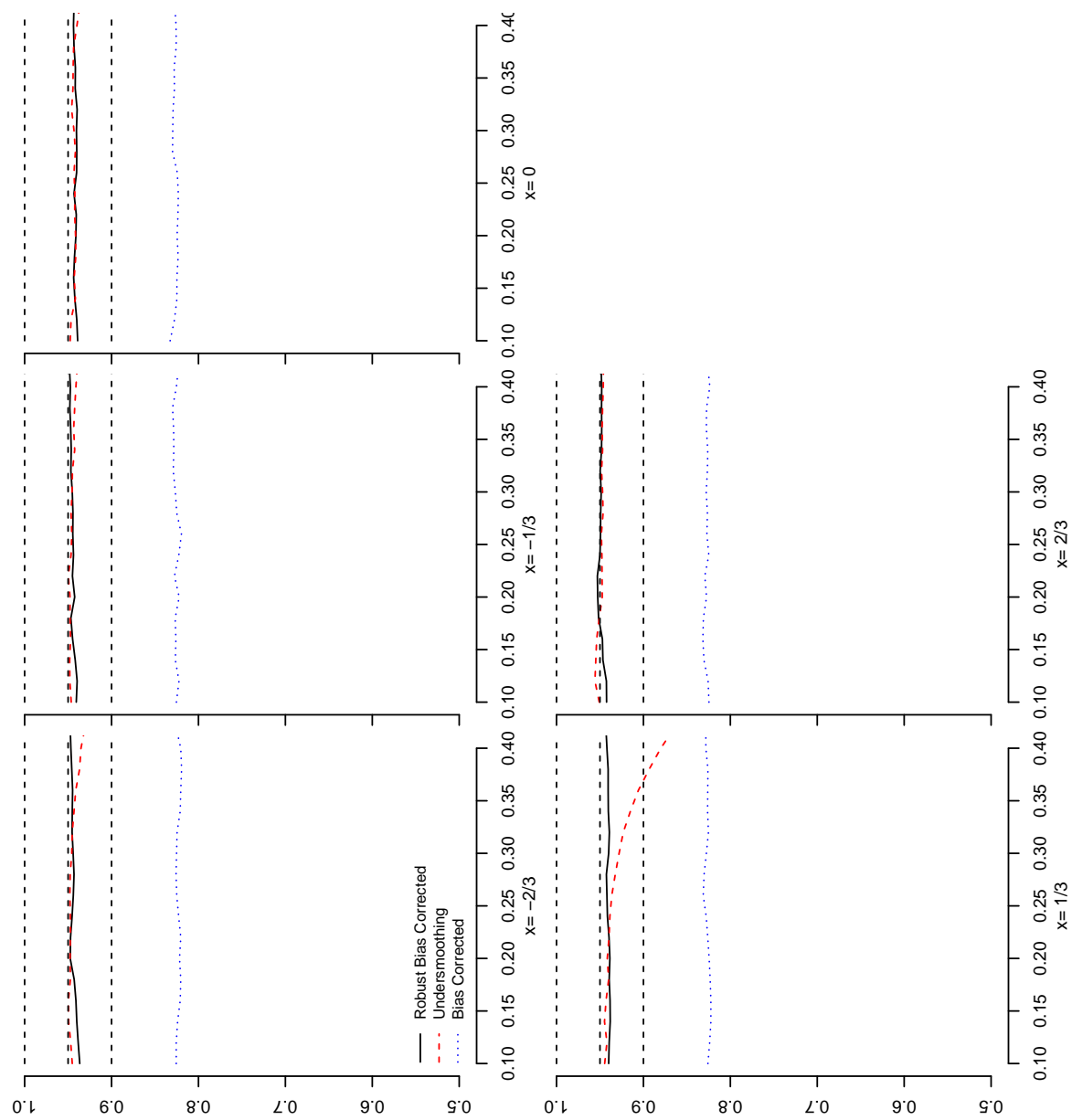


Figure S.II.8: Average Interval Length of 95% Confidence Intervals - Model 1

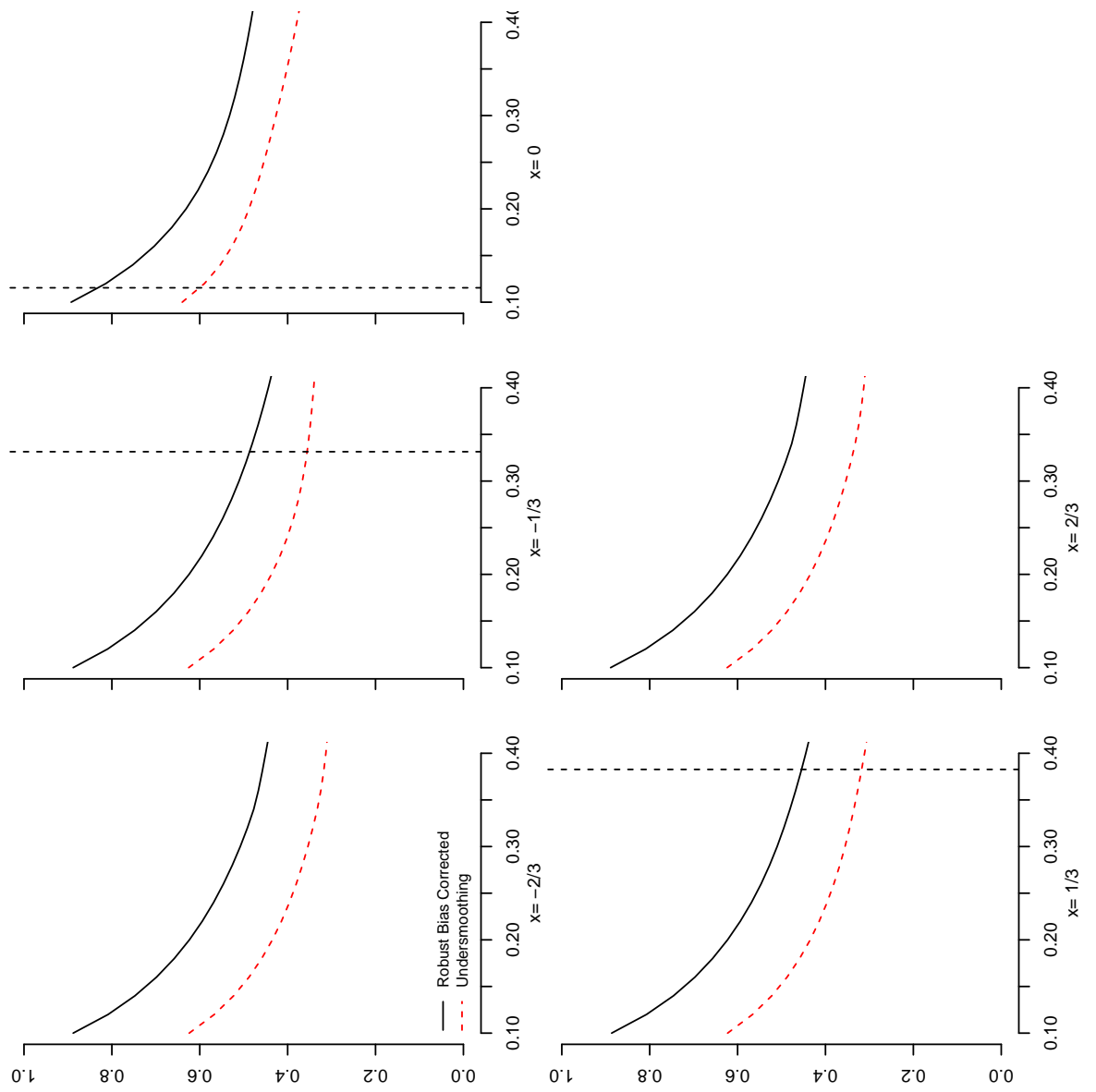




Figure S.II.9: Average Interval Length of 95% Confidence Intervals - Model 2

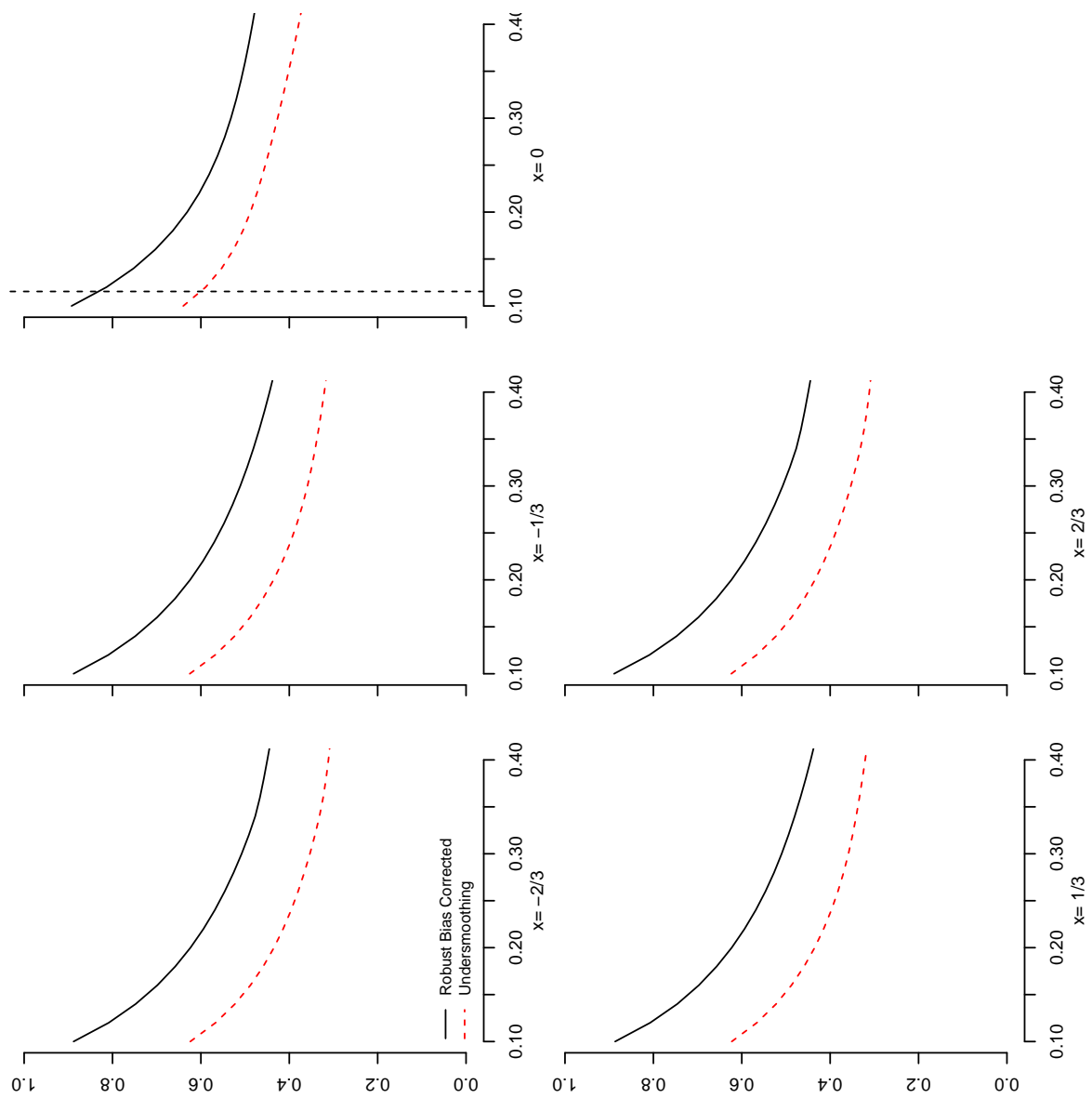


Figure S.II.10: Average Interval Length of 95% Confidence Intervals - Model 3

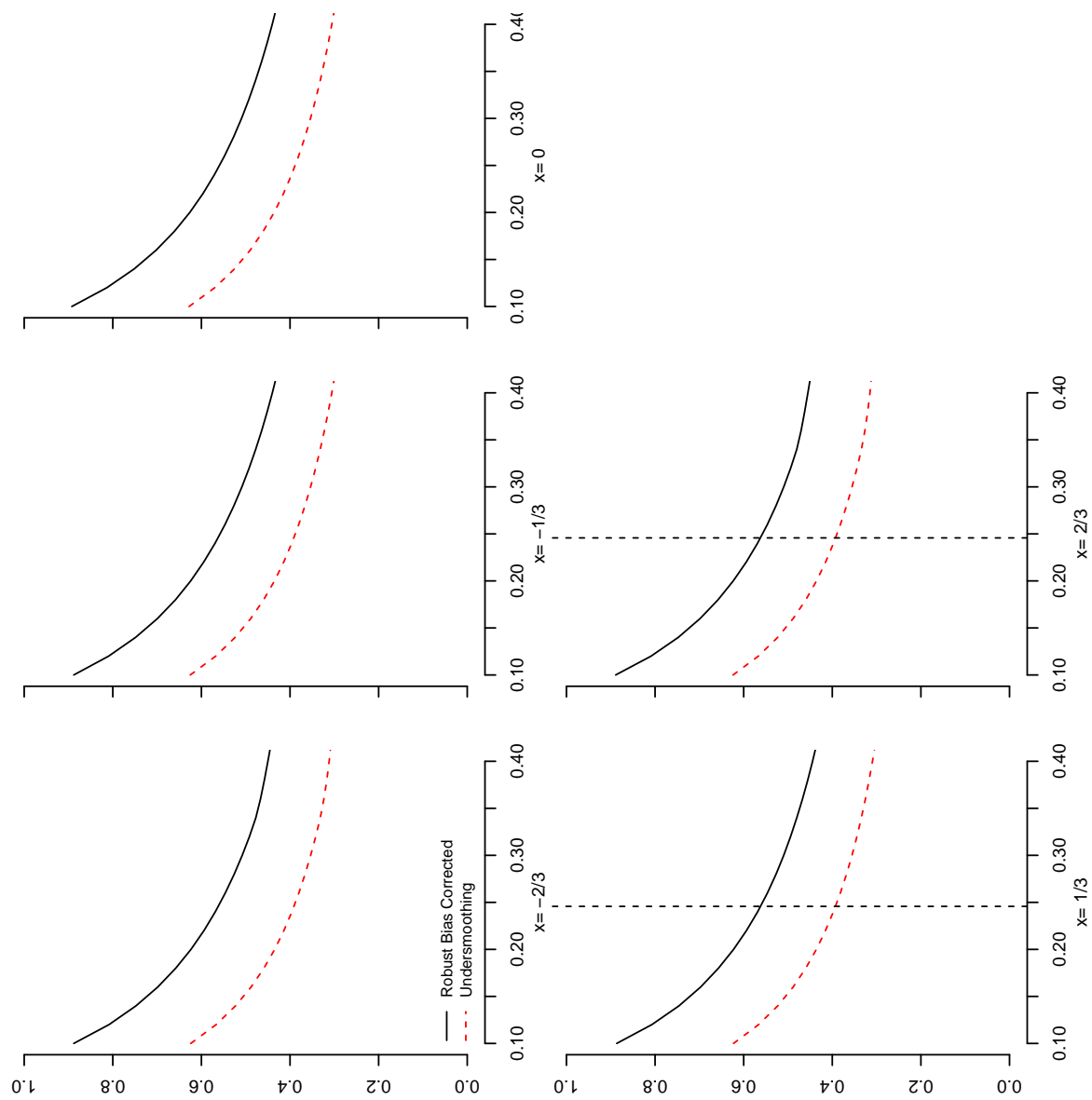


Figure S.II.11: Average Interval Length of 95% Confidence Intervals - Model 4

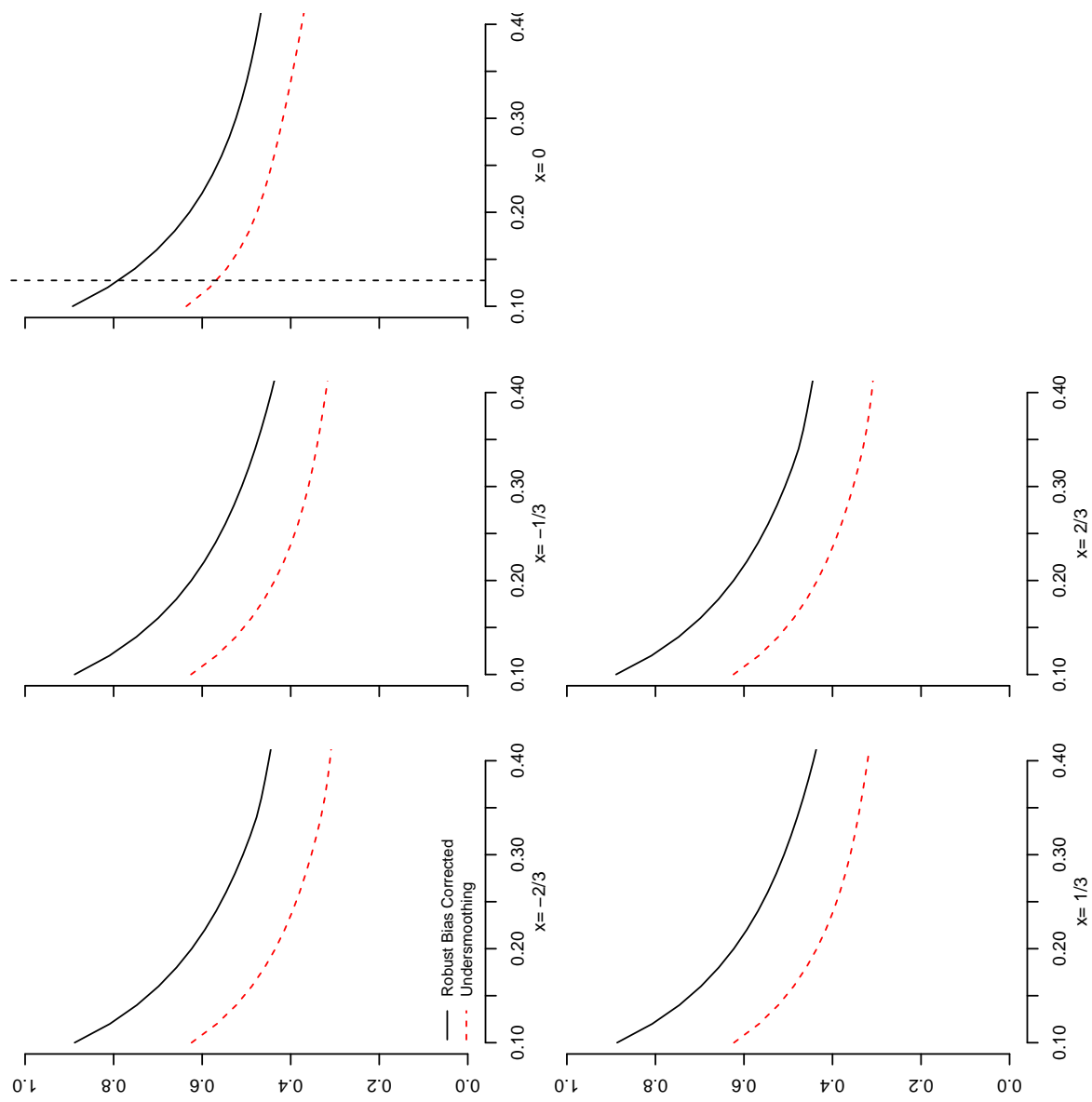


Figure S.II.12: Average Interval Length of 95% Confidence Intervals - Model 5

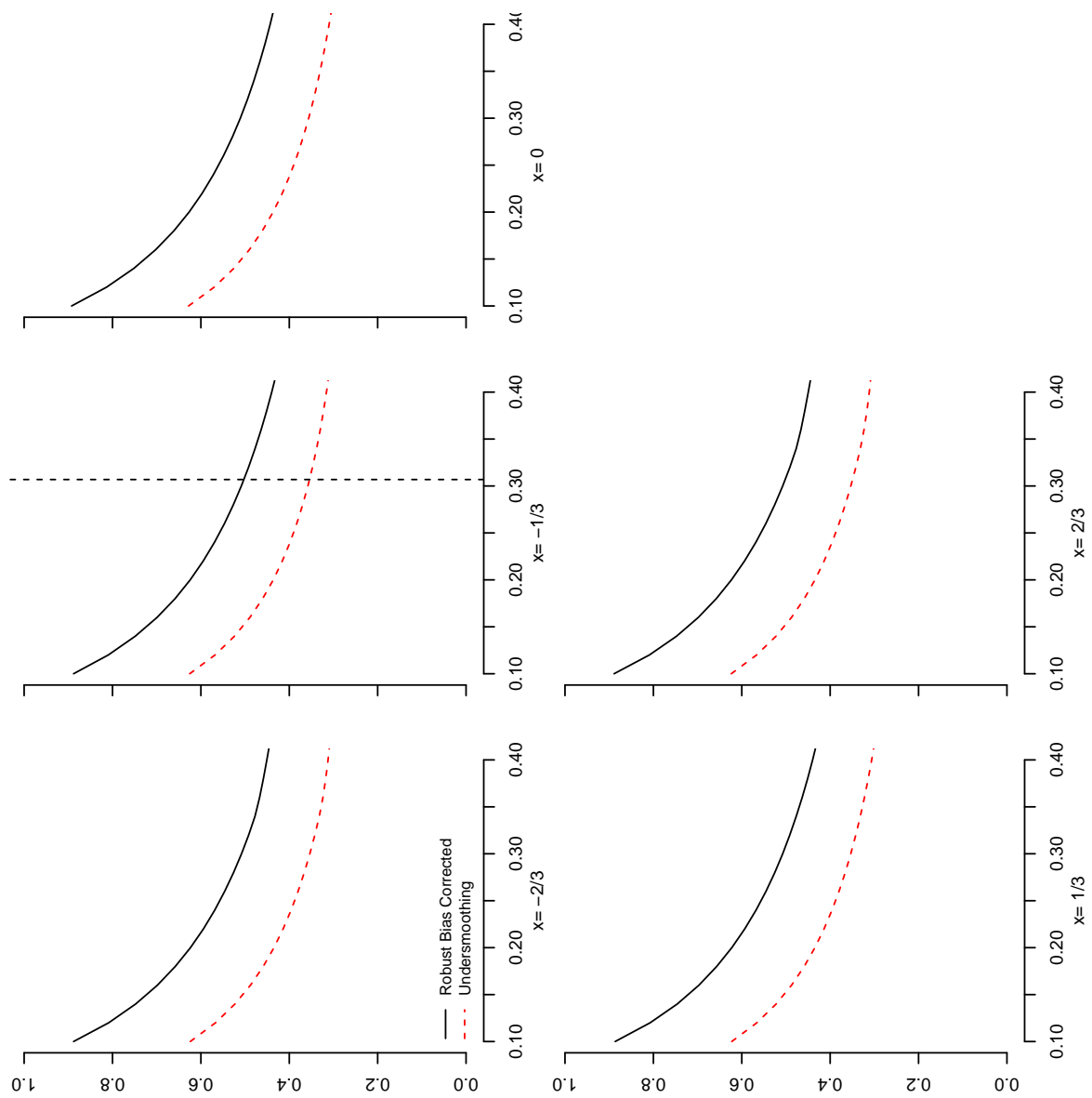


Figure S.II.13: Average Interval Length of 95% Confidence Intervals - Model 6

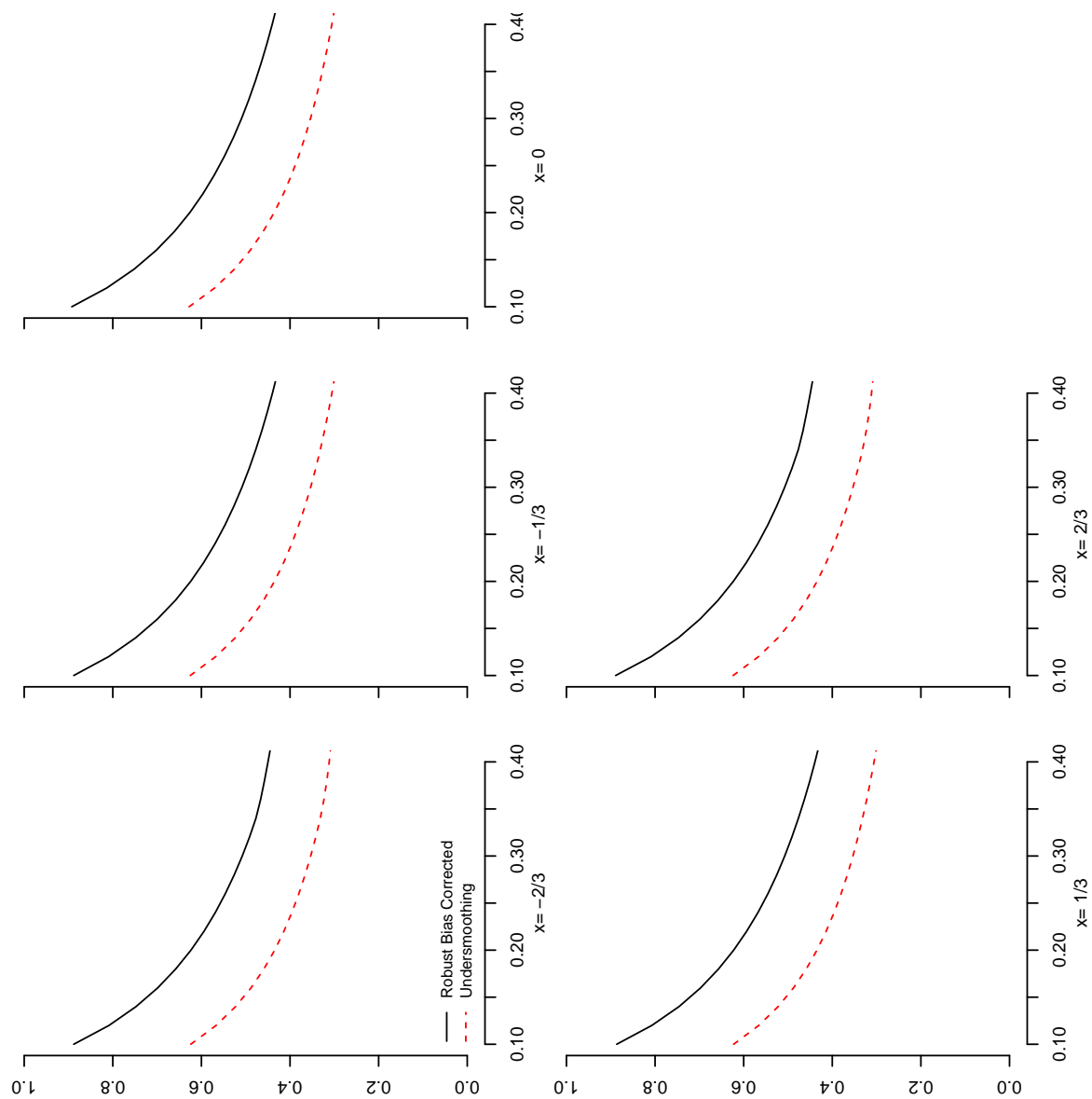


Figure S.II.14: Empirical Coverage and Average Interval Length of 95% Confidence Intervals - Model 1

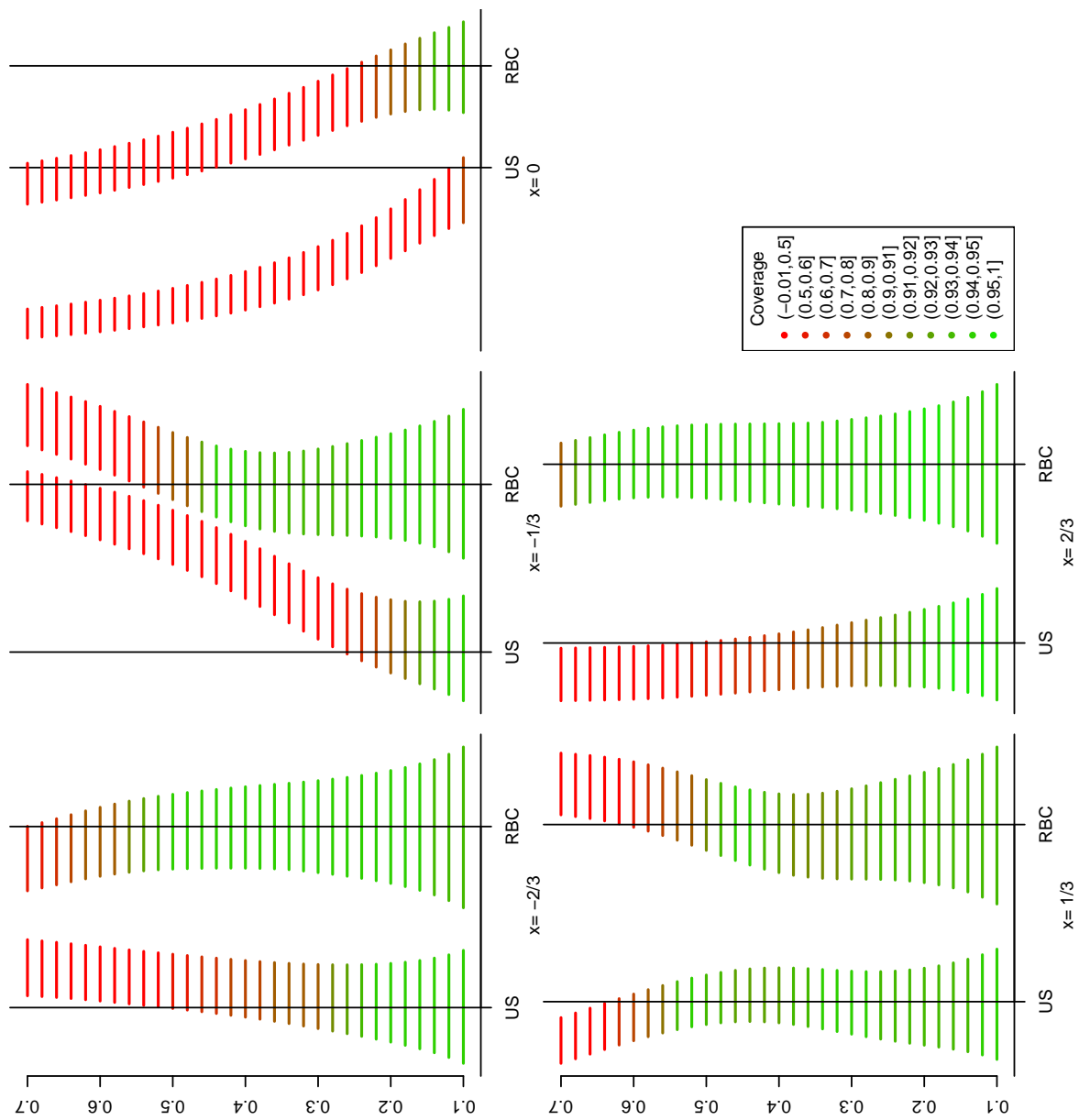


Figure S.II.15: Empirical Coverage and Average Interval Length of 95% Confidence Intervals - Model 2

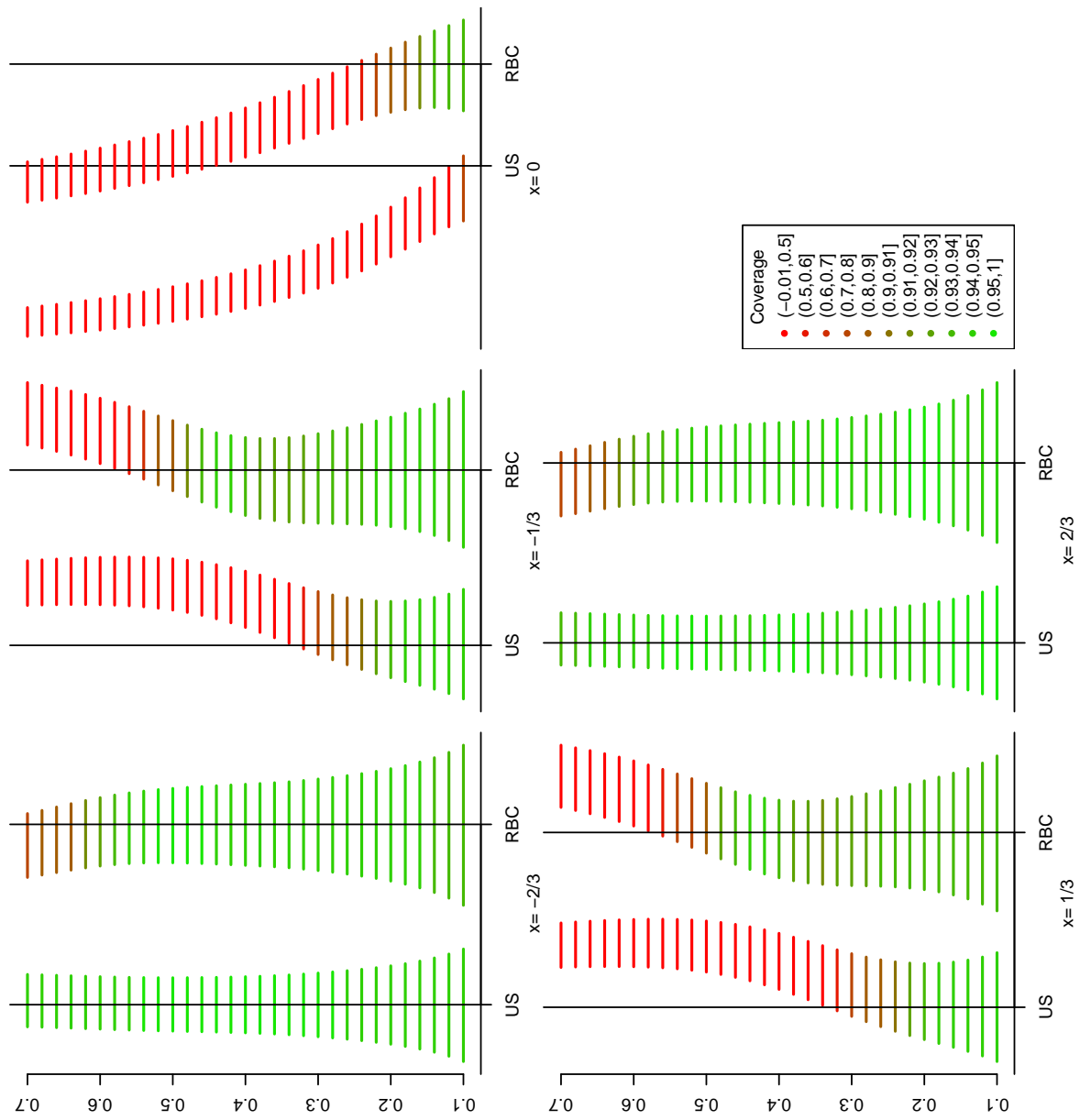


Figure S.II.16: Empirical Coverage and Average Interval Length of 95% Confidence Intervals - Model 3

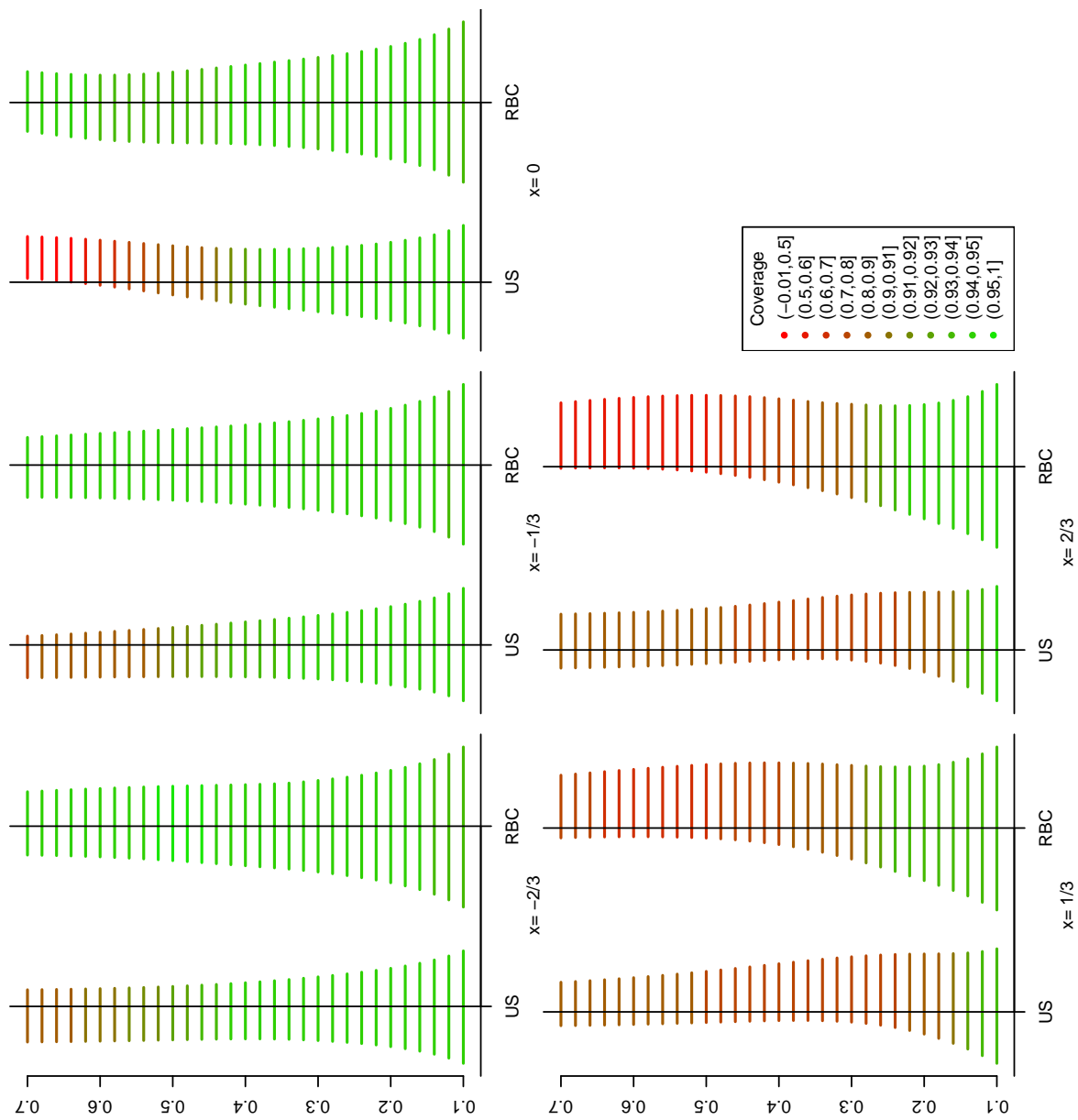




Figure S.II.17: Empirical Coverage and Average Interval Length of 95% Confidence Intervals - Model 4

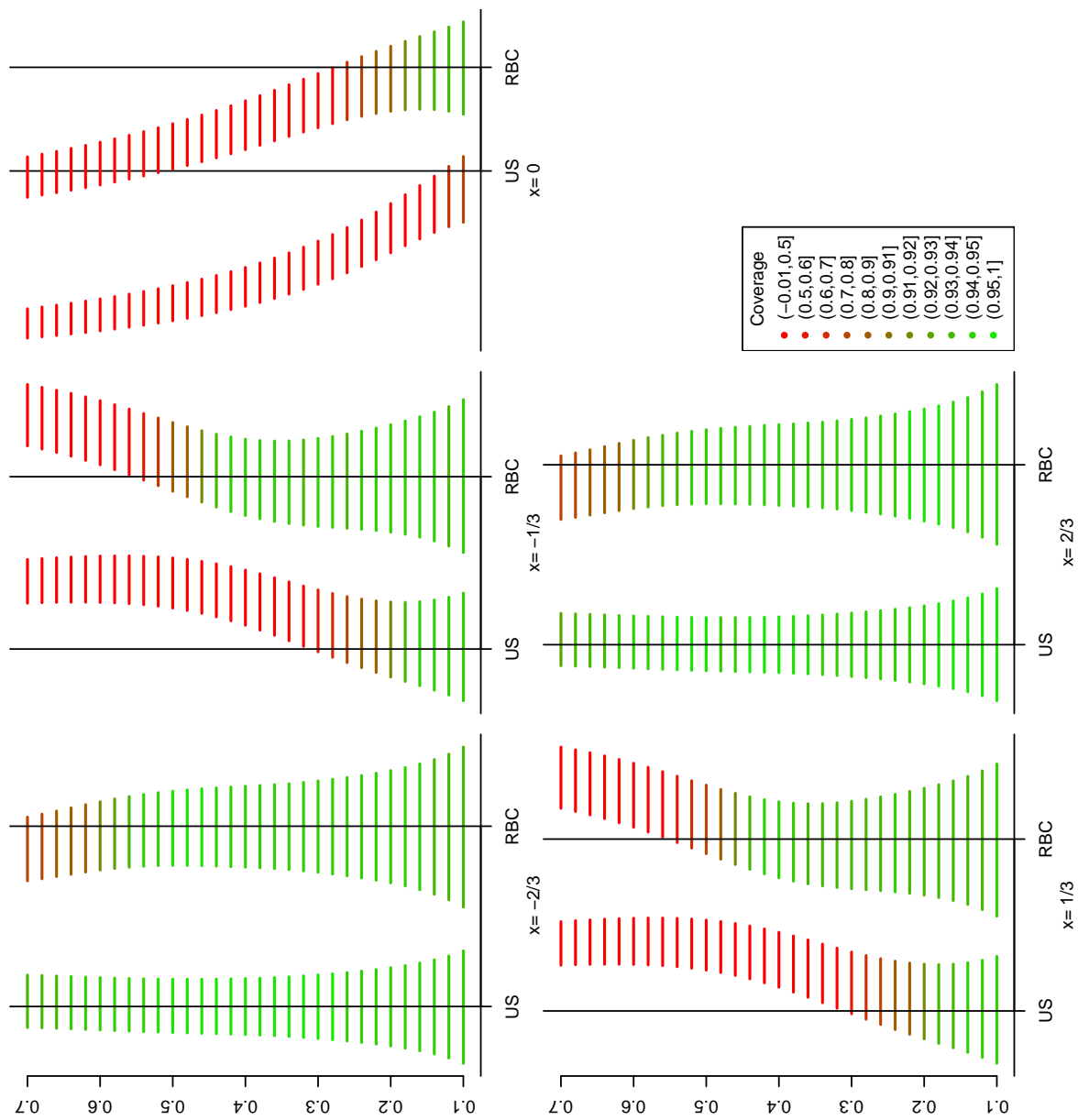


Figure S.II.18: Empirical Coverage and Average Interval Length of 95% Confidence Intervals - Model 5

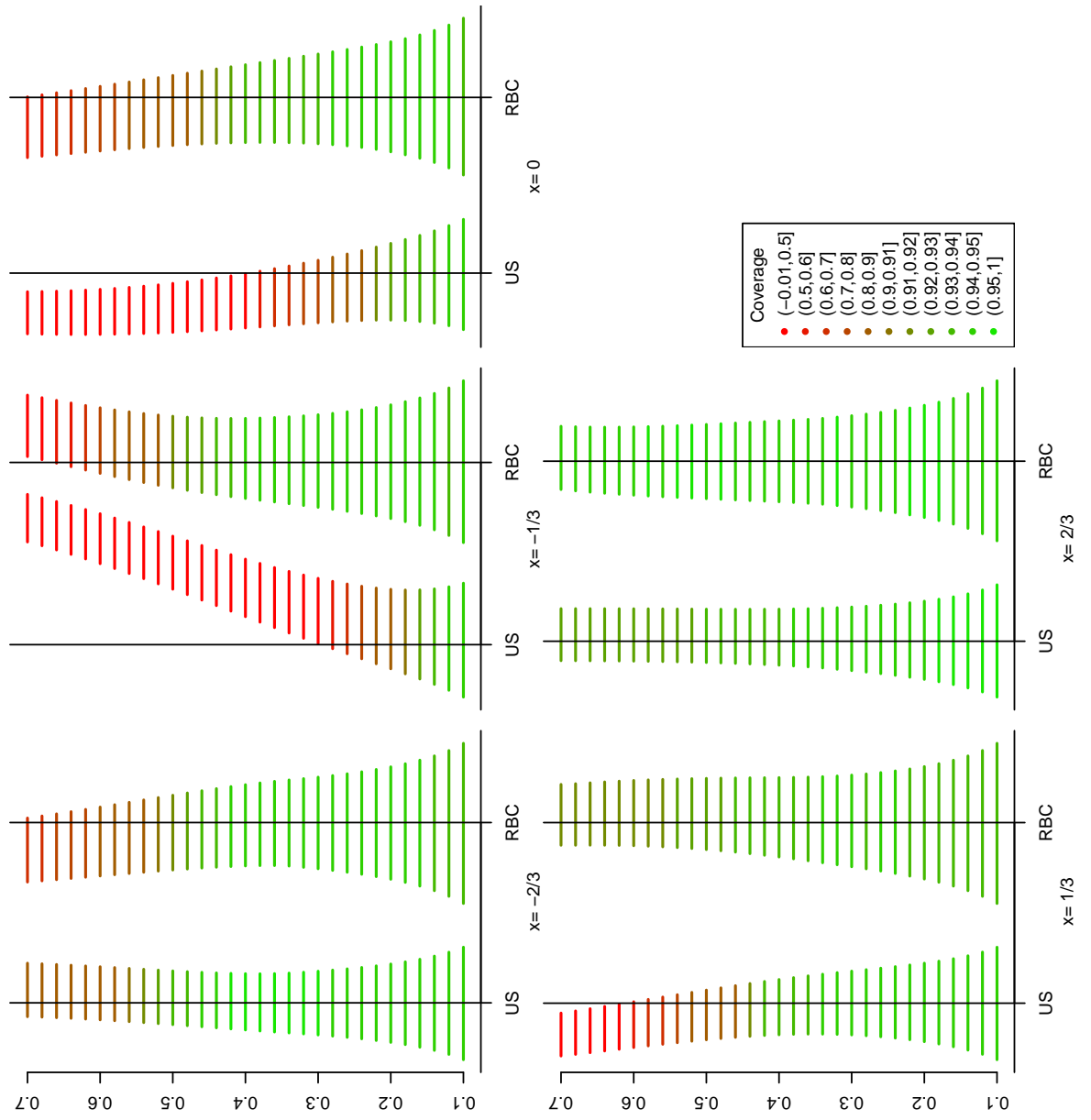


Figure S.II.19: Empirical Coverage and Average Interval Length of 95% Confidence Intervals - Model 6

